

EXHIBIT A

Readers Prefer Outputs of AI Trained on Copyrighted Books over Expert Human Writers

Tuhin Chakrabarty¹, Jane C. Ginsburg², Paramveer Dhillon^{3,4}

¹Department of Computer Science and AI Innovation Institute, Stony Brook University.

²Columbia Law School.

³School of Information Science, University of Michigan.

⁴MIT Initiative on the Digital Economy.

Corresponding authors: tchakrabarty@cs.stonybrook.edu, ginsburg@law.columbia.edu, dhillionp@umich.edu

The use of copyrighted books for training AI models has led to numerous lawsuits from authors concerned about AI’s ability to generate derivative content. Yet it’s unclear whether these models can generate high quality literary text while emulating authors’ styles/voices. To answer this we conducted a preregistered study comparing MFA-trained expert writers with three frontier AI models: ChatGPT, Claude, and Gemini in writing up to 450 word excerpts emulating 50 award-winning authors’ (including Nobel laureates, Booker Prize winners, and young emerging National Book Award finalists) diverse styles. In blind pairwise evaluations by 159 representative expert (MFA candidates from top U.S. writing programs) and lay readers (recruited via Prolific), AI-generated text from in-context prompting was strongly disfavored by experts for both stylistic fidelity (odds ratio [OR]=0.16, $p < 10^{-8}$) and writing quality (OR=0.13, $p < 10^{-7}$) but showed mixed results with lay readers. However, fine-tuning ChatGPT on individual author’s complete works completely reversed these findings: experts now favored AI-generated text for stylistic fidelity (OR=8.16, $p < 10^{-13}$) and writing quality (OR=1.87, $p=0.010$), with lay readers showing similar shifts. These effects are robust under cluster-robust inference and generalize across authors and styles in author-level heterogeneity analyses. The fine-tuned outputs were rarely flagged as AI-generated (3% rate versus 97% for in-context prompting) by state-of-the-art AI detectors. Mediation analysis reveals this reversal occurs because fine-tuning eliminates detectable AI stylistic quirks (e.g., cliché density) that penalize in-context outputs, altering the relationship between AI detectability and reader preference. While we do not account for additional costs of human effort required to transform raw AI output into cohesive, publishable prose, the median fine-tuning and inference cost of \$81 per author represents a dramatic 99.7% reduction compared to typical professional writer compensation. Author-specific fine-tuning thus enables non-verbatim AI writing that readers prefer to expert human writing, thereby providing empirical evidence directly relevant to copyright’s fourth fair-use factor, the “effect upon the potential market or value” of the source works.

Keywords: Generative AI, Copyright Law, Fair Use, Future of Work, AI Detection, AI and Society, Behavioral Science, Labor Market Impact

1 Introduction

The U.S. publishing industry supports hundreds of thousands of jobs while generating \$30 billion in yearly revenue, contributing to the larger American copyright sectors that account for \$2.09 trillion in annual GDP contributions (1). Adult Fiction and nonfiction books alone accounted for \$6.14 billion in 2024 (2). This economically important sector now faces an unprecedented challenge: its core products have become essential training data for generative-AI systems. Models trained on well-edited books produce more coherent, accurate responses—something crucial to creating the

illusion of intelligence (3). Most technology companies building AI use massive datasets of books, typically without permission or licensing (4), and frequently from illegal sources. In the recent copyright lawsuit *Bartz v. Anthropic* (5), Judge Alsup noted that Anthropic acquired at least five million books from LibGen and two million from Pirate Library Mirror (PiLiMi). Anthropic also used Books3¹—a dataset of approximately 191,000 books also used by Meta and Bloomberg to train their language models.² This unauthorized use has sparked outrage among authors (6), triggering dozens of lawsuits against technology companies including OpenAI, Anthropic, Microsoft, Google, and Meta.

Generative AI systems such as ChatGPT that can be prompted to create new text at scale are qualitatively unlike most historical examples of automation technologies (7). They can now solve Olympiad-level Geometry (8), achieve an impressive rating of 2700 on Codeforces, one of the most challenging coding competition platforms (9), and deliver medical guidance that meets professional healthcare standards (10, 11). Recent findings from Microsoft's Occupational Implications of Generative AI (12), Anthropic's Economic Index (13) and OpenAI (14) reveal AI usage primarily concentrates in writing tasks. This concentration threatens creative writing professionals in particular—novelists, poets, screenwriters, and content creators who shape cultural narratives and human expression. Based on U.S. Bureau of Labor Statistics May 2023 national estimates, creative writing constitutes almost 50% of writing jobs (15), making these positions especially vulnerable to GenAI-based automation, as the writing community has warned (16).

While it is widely established that most frontier large language models (LLMs) have been trained on copyrighted books, it remains unclear whether such training can produce expert-level creative writing. Past research has shown that AI cannot produce highbrow literary fiction or creative nonfiction through prompting alone when compared to professionally trained writers (17). More recent work from (18) demonstrates that AI-generated creative writing still remains characterized by clichés, purple prose, and unnecessary exposition. Additionally, relying on GenAI for creative writing reduces the collective diversity of novel content (19). AI often produces formulaic, mediocre creative writing because it lacks the distinctive personal voice that typically distinguishes one author from another (20). As Pulitzer Fiction finalist Vauhini Vara observes, “ChatGPT’s voice is polite, predictable, inoffensive, upbeat. Great characters, on the other hand, aren’t polite; great plots aren’t predictable; great style isn’t inoffensive; and great endings aren’t upbeat” (21). To address this limitation, practitioners now increasingly prompt AI systems to perform style/voice mimicry by emulating specific writers’ choices (22). This practice has become so common that a fantasy author recently published a novel containing an accidentally included AI prompt requesting emulation of another writer’s style (23). While the effectiveness of such stylistic emulation remains contested, the more pressing question concerns whether style/voice mimicry genuinely improves AI-generated text quality and whether readers—both experts and non-experts—perceive these improvements as meaningful.

To address this question, we conducted a preregistered behavioral study comparing MFA-trained expert writers with frontier large language models. Historically, top MFA programs have produced many prizewinning American writers (24). Eminent literary agent Gail Hochman of Brandt & Hochman said, “We look favorably on anyone who has an M.F.A., simply because it shows they’re serious about their writing” (25). This elite sample of MFA trained expert writers provides a conservative test—if AI can compete with the best emerging talent, the disruption to average writers is likely even greater. We selected closed-source LLMs readily accessible to users without technical expertise: GPT-4o, Claude 3.5 Sonnet, and Gemini 1.5 Pro.³ Additionally our results would hold across different LLM versions in the future as LLM outputs have been shown to be homogeneous and have not improved over time (26, 27, 28, 29). Both human experts and LLMs were given the same task: write an excerpt of up to 450 words emulating the style and voice of one of 50 internationally acclaimed authors, including Nobel laureates, Booker Prize winners, and Pulitzer Prize winners, spanning multiple continents and cultures.⁴ We tested two AI conditions: (1) in-context prompting, where models received the same instructions as human experts, and (2) fine-tuning, where models were additionally trained on each author’s complete oeuvre.⁵ Expert and lay readers performed blind pairwise evaluations (30, 31, 32) of <Human-AI> excerpts on writing quality and stylistic fidelity. This design addresses three preregistered research questions: (1) Can AI match expert performance in writing quality and stylistic fidelity across both conditions? (2) Do expert and lay readers show similar preference patterns? (3) Does AI detectability correlate with human quality

¹Now removed after a legal complaint by anti-piracy group, the Rights Alliance.

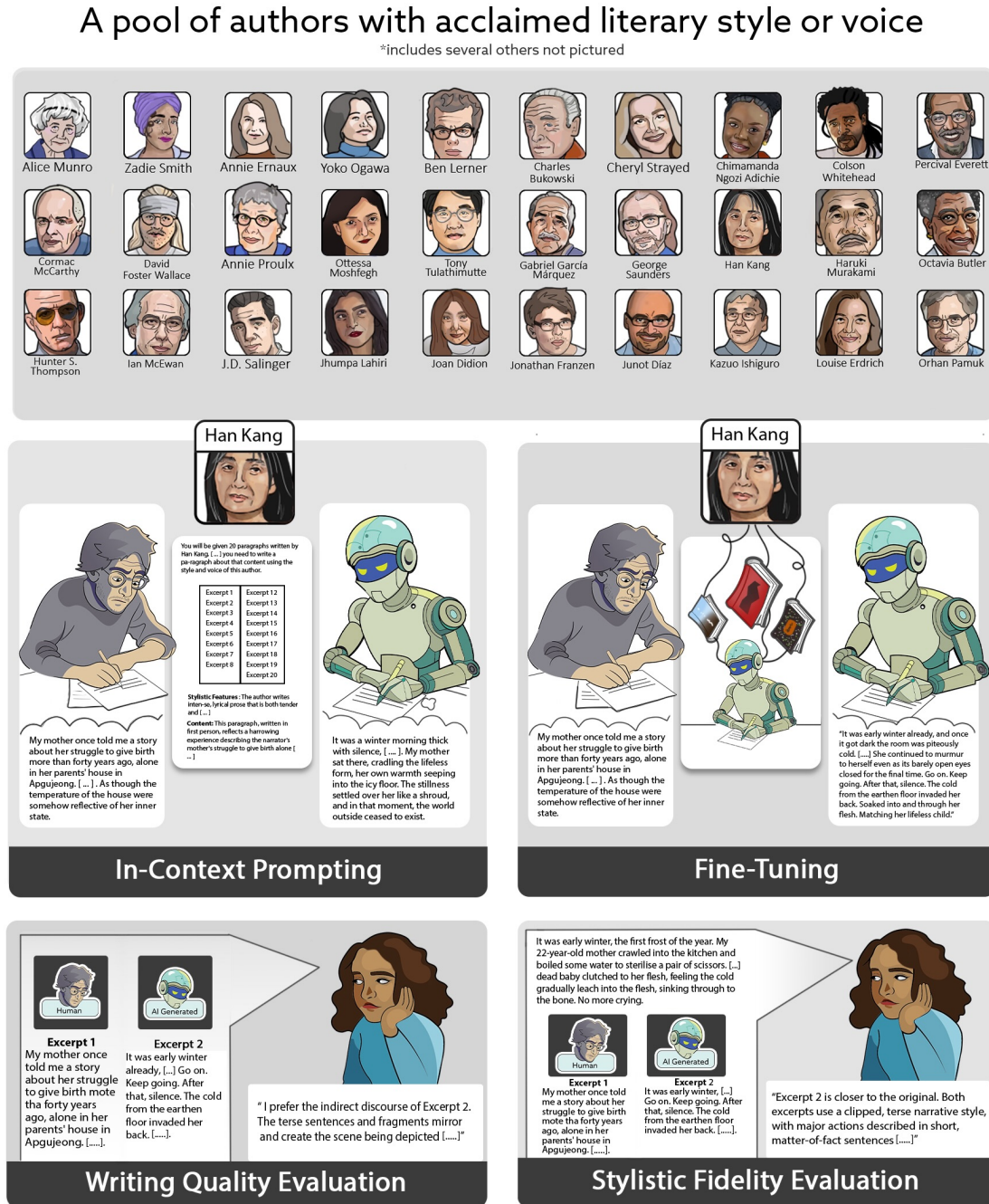
²The *Bartz v. Anthropic* lawsuit also revealed how Anthropic cut millions of print books from their bindings, scanned them into digital files, and threw away the originals solely for the purpose of training Claude.

³We also tried open-weight Llama 3.1 model but its empirical performance was not good at following long context instructions at the time of our study.

⁴Among others, our study includes Nobel laureates Han Kang and Annie Ernaux; Booker Prize winners Salman Rushdie, Margaret Atwood, and George Saunders; and Pulitzer Prize winners Junot Diaz and Marilynne Robinson.

⁵For authors writing in non-English languages (Han Kang, Yoko Ogawa, Annie Ernaux, Haruki Murakami), we used the same translator’s work across all books to maintain voice consistency.

judgments, and does fine-tuning remove this correlation? Our full experimental setup is shown in Figure 1.



Preference Judgements from Expert and Lay Readers

Figure 1: Figure showing our study design. (1) Select a target author and prompt. (2) Generate upto 450-word candidate excerpts from MFA experts and from LLMs under two settings: in-context prompting (instructions + few-shot examples) and author-specific fine-tuning (model fine-tuned on that author's works). (3) Readers (experts and lay) perform blinded, pairwise forced-choice evaluations on two outcomes: stylistic fidelity to the target author and overall writing quality. Pair order and left/right placement are randomized on every trial.

2 Results

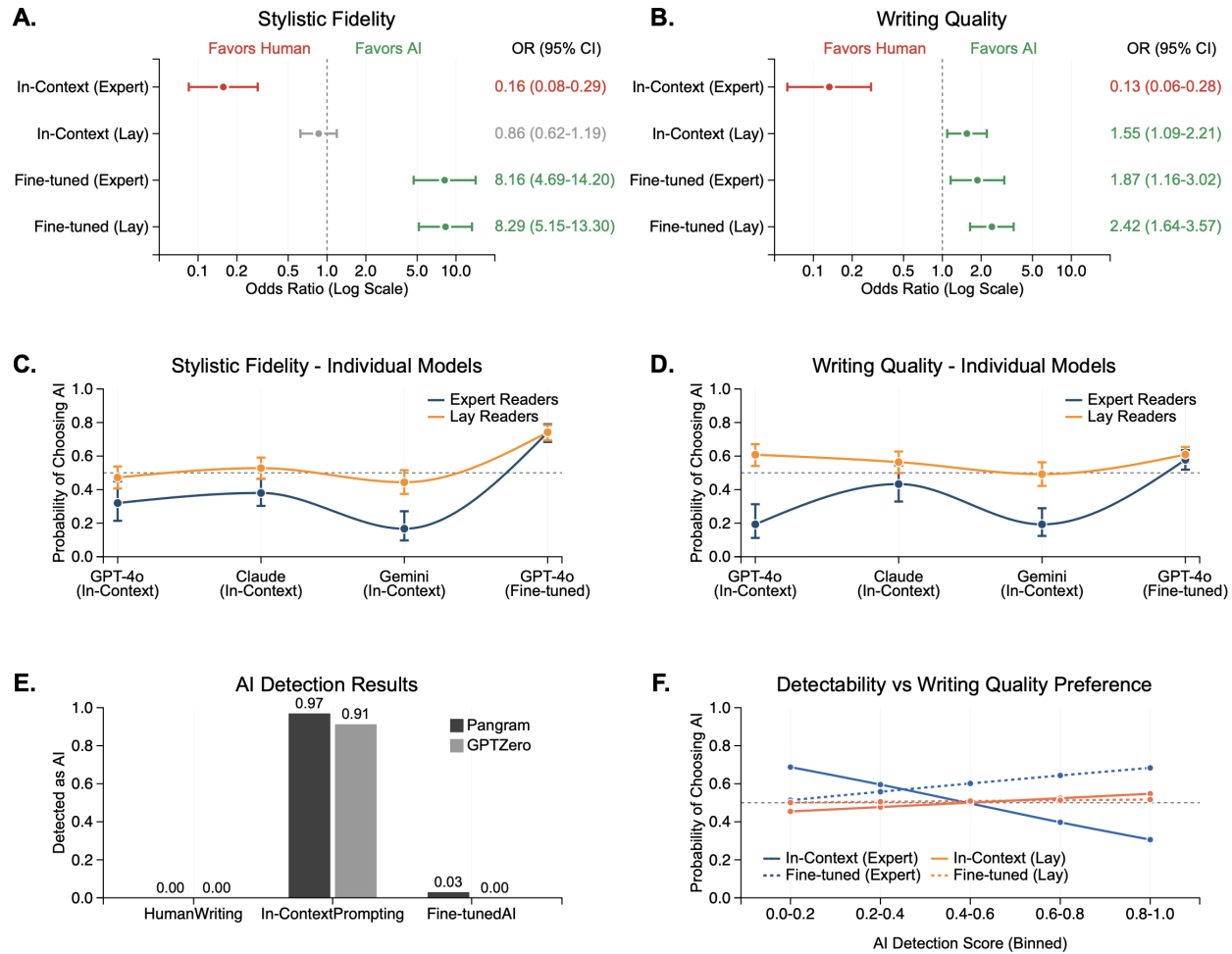


Figure 2: (A-B) Forest plots showing odds ratios (OR) and 95% confidence intervals comparing AI and human experts in pairwise evaluations of stylistic fidelity (A) and writing quality (B) where values >1 favor AI and values <1 favor humans. Expert readers show preference for human writing when prompted in an in-context setting (OR = 0.16 and 0.13) but that changes when AI is fine-tuned (OR = 8.16 and 1.87). Lay readers have a harder time discriminating, given how they prefer the quality of AI writing even for in-context prompting (OR = 1.55). (C-D) Probability of choosing AI excerpts across individual language models for stylistic fidelity (C) and writing quality (D). Error bars represent 95% confidence intervals. Dashed line indicates chance performance (50%). (E) AI detection accuracy with chosen threshold of $\tau=0.9$ using two state-of-the-art AI detectors (Pangram and GPTZero). Human written text was never misclassified (0.00), in-context AI was detected with 97% accuracy by Pangram and 91% by GPTZero, but fine-tuned AI evaded detection 97% of the time (0.03) for Pangram and 100% of the time in case of GPTZero. (F) Relationship between AI detectability (Pangram) and preference for writing quality across detection score bins. For in-context prompting setup, higher detection scores correlated with lower AI preference (negative slope). This relationship disappeared when AI is fine-tuned (flat slopes). Fine-tuning on an author's complete oeuvre eliminates stylistic "AI" quirks while achieving expert-level performance. $n = 28$ expert readers, 131 lay readers; 3,840 pairwise comparisons with robust clustered standard errors.

2.1 Overall Performance Comparisons

Our final data consist of 3,840 paired-choice tasks, with judgments from 28 MFA experts and 131 lay readers. We fit a logit model for each outcome and condition and include dummies for writer type, reader group, and their interaction. We further employ CR2 cluster-robust standard errors (33) clustered at the reader-level to account for within-reader

correlation in ratings. Our hypotheses, outcomes, design, and analysis closely follow our OSF pre-registration (SI Sections S4-S8); deviations are detailed in SI Section S9.

Figure 2A–B presents the odds ratios, with corresponding predicted probabilities shown in Figure 2C–D. Under in-context prompting, expert readers demonstrated strong preference for human-written text. Odds ratios were 0.16 (95% CI: 0.08–0.29, $p < 10^{-8}$) for stylistic fidelity and 0.13 (95% CI: 0.06–0.28, $p < 10^{-7}$) for writing quality, indicating six- to eight-fold preferences for human excerpts (Fig. 2A–B). Lay readers showed no significant preference regarding stylistic fidelity (OR = 0.86, 95% CI: 0.62–1.19, $p = 0.37$) but favored AI-generated text for writing quality (OR = 1.55, 95% CI: 1.09–2.21, $p = 0.014$), selecting AI excerpts in 61% of writing quality trials (Fig. 2C–D). Inter-rater agreement reflected this divergence: expert readers achieved $\kappa = 0.58$ for stylistic fidelity and $\kappa = 0.41$ for writing quality, while lay readers showed minimal agreement among themselves ($\kappa = 0.12$ and $\kappa = 0.15$, respectively).⁶ The writer-type \times reader-type interaction was significant for both outcomes ($\chi^2_{(3)} = 24.9$, $p = 1.6 \times 10^{-5}$ for fidelity; $\chi^2_{(3)} = 37.6$, $p = 3.5 \times 10^{-8}$ for quality).

Fine-tuning on authors' complete works reversed these preferences. For expert readers, the odds of selecting the AI excerpt were 8.16 times the odds of selecting the human excerpt for stylistic fidelity (OR = 8.16, 95% CI: 4.69–14.2, $p < 10^{-13}$) and 1.87 times the odds for writing quality (OR = 1.87, 95% CI: 1.16–3.02, $p = 0.010$). Lay readers showed comparable shifts in their preferences (stylistic fidelity OR = 8.29, 95% CI: 5.15–13.3, $p < 10^{-17}$; writing quality OR = 2.42, 95% CI: 1.64–3.57, $p < 10^{-5}$). Model-based predicted AI win probabilities converged across groups to about 0.74 for stylistic fidelity and 0.58–0.61 for writing quality (Fig. 2C–D), and the writer-type \times reader-type interaction was no longer significant in the fine-tuned models (both $\chi^2_{(1)} < 1$, $p > 0.40$). Inter-rater agreement among experts increased ($\kappa = 0.67$ for writing quality; $\kappa = 0.54$ for stylistic fidelity), while agreement among lay readers remained low ($\kappa = 0.07$ and $\kappa = 0.22$, respectively).

2.2 AI Detection and Stylometric Analysis

We probe whether differences in AI detectability can account for these preference reversals. Pangram, a state-of-the-art AI detection tool (34, 35, 36), correctly classified 97% of in-context prompted texts as machine-generated but only 3% of fine-tuned texts were classified as AI-generated (Fig. 2E).⁷

Higher AI-detection scores strongly predicted lower preference rates among expert readers in the in-context prompting condition. For stylistic fidelity, each unit increase in detection score reduced the odds of selecting an excerpt by a factor of 6.3 ($\beta = -1.85 \pm 0.29$, $p < 10^{-9}$); a similar pattern held for writing quality ($\beta = -2.01 \pm 0.33$, $p < 10^{-9}$). Fine-tuning largely eliminated this negative relationship between detectability and preference (Pangram \times Setting for style: $\beta = 2.56 \pm 0.81$, $p = 0.002$; for quality: $\beta = 2.90 \pm 0.88$, $p < 0.001$). Two-stage mediation analysis (Fig. 4A) demonstrated that stylometric features, particularly cliché density (See Section S3.2 in SI), mediated 16.4% of the detection effect on preference before fine-tuning but a statistically insignificant 1.3% (95% CI includes zero) afterward, indicating that fine-tuning eliminates rather than masks artificial stylistic signatures.

2.3 Author-Level Performance Heterogeneity

Next, we disaggregated our data to unpack heterogeneity at the level of individual authors. Of 30 fine-tuned author models, 27 exceeded parity for stylistic fidelity (median win rate = 0.74, IQR: 0.63–0.86) and 23 for writing quality (median = 0.58, IQR: 0.54–0.74). Using 95% Jeffreys intervals, fine-tuned models significantly outperformed human writers for 19 authors on stylistic fidelity and 10 authors on writing quality (Fig. 3A). These performance differences showed no systematic relationship with fine-tuning corpus size (Fig. 3B). The fine-tuning premium, i.e., the increase in AI preference rates relative to in-context prompting, ranged from –13.7 to +70.8 percentage points for stylistic fidelity (29 of 30 positive) and from –20.8 to +50.0 percentage points for writing quality (22 of 30 positive). This “premium” likewise showed no correlation with fine-tuning corpus size (Pearson $r < 0.1$ for both outcomes; Fig. 4B).

⁶Given that lay assessments of literary quality and style reflect inherently diverse tastes, we anticipated lower inter-rater agreement compared to expert readers.

⁷GPTZero, another state-of-the-art AI detection tool showed comparable performance but had higher false-positive-rate, so for our subsequent analyses we stuck to Pangram.

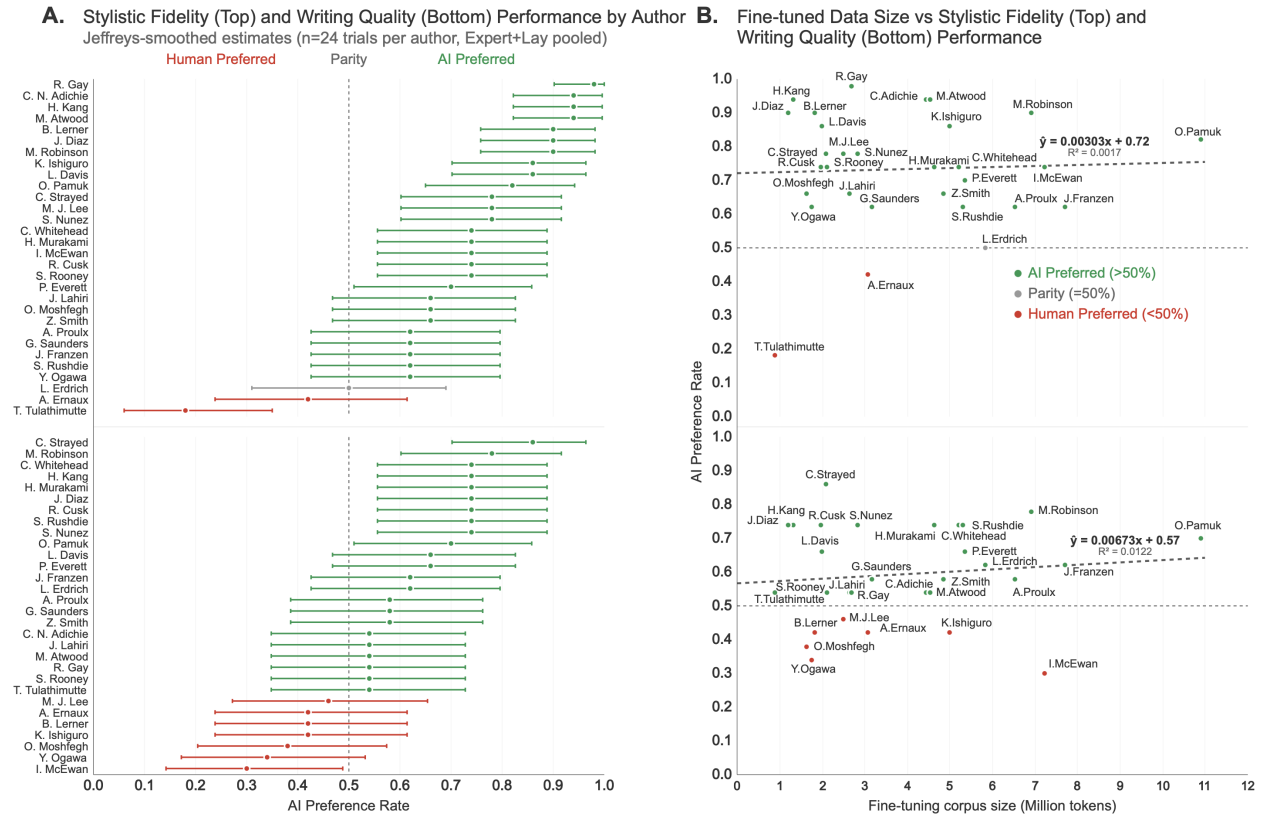


Figure 3: Author-level AI preference and its association with fine-tuning corpus size (fidelity and quality) (A) For each fine-tuned author, the share of blinded pairwise trials in which the AI excerpt was preferred over the human (MFA expert) on stylistic fidelity (top) and overall quality (bottom). Points show Jeffreys-prior estimates ($(k + 0.5)/(n + 1)$); vertical bars are 95% Jeffreys intervals ($\text{Beta}(\frac{1}{2}, \frac{1}{2})$); the dotted line at 0.5 marks human–AI parity. Readers are pooled (experts and lay). (B) AI preference rate versus the fine-tuning corpus size for that author (million tokens), shown for stylistic fidelity (top) and overall quality (bottom). Each point is a fine-tuned author; the line is an OLS fit with heteroskedasticity-robust standard errors (no CI displayed). Slopes are near zero in both panels, indicating little association between corpus size (in this range) and AI preference.

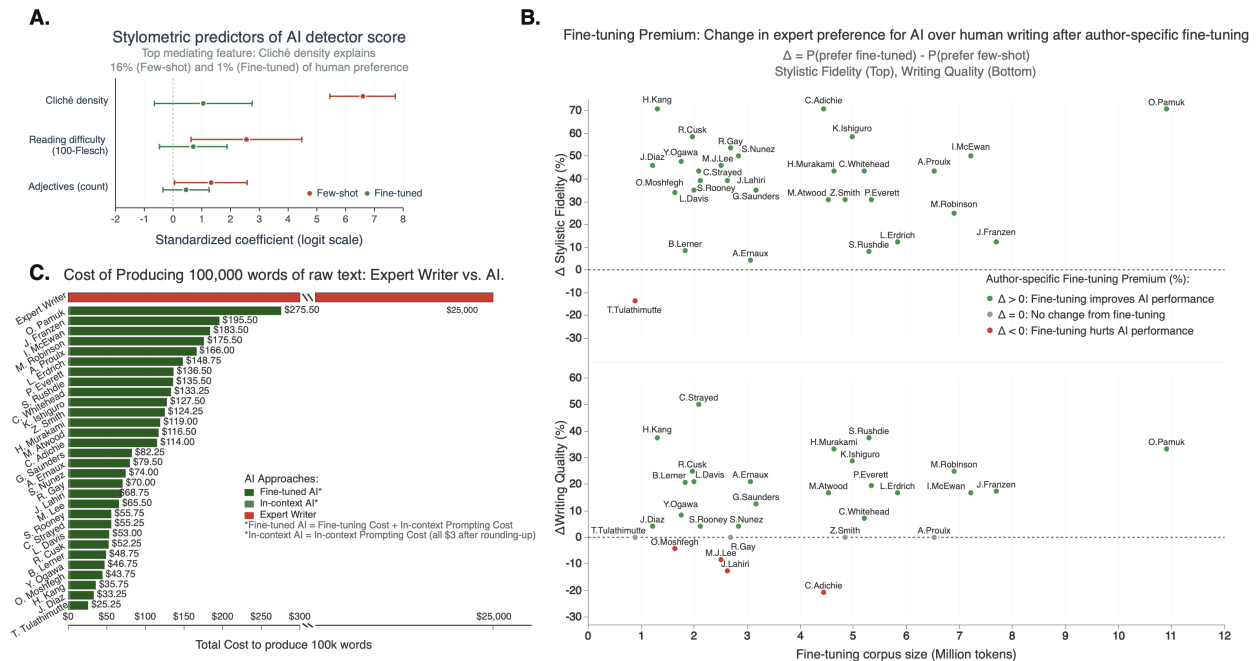


Figure 4: Fine-tuning substantially reduces stylistometric signatures of AI text, improves stylistic fidelity and perceived writing quality over in-context prompting, and substantially cuts costs of producing first draft versus professional writers. (A) Mediation analysis linking stylistometric features \rightarrow AI-detector score \rightarrow human preference. Standardized logistic coefficients with 95% CIs are shown for three features for in-context prompting (red) and fine-tuned models (green). Cliché density mediates 16.4% of the detector effect on choice for in-context prompting but only 1.3% for author fine-tuned models; all three features together mediate 25.4% vs -3.2% for in-context prompted and fine-tuned models respectively. (B) “Fine-tuning premium,” defined as $\Delta = P(\text{prefer fine-tuned over human}) - P(\text{prefer in-context over human})$, as a function of fine-tuning corpus size. Top: stylistic fidelity; bottom: writing quality. Points are authors; colors denote improvement (green, $\Delta \geq 0$), no change (gray), or degradation (red, $\Delta < 0$). Median Δ : +41.7 (fidelity) and +16.7 percentage points (quality). (C) Cost to produce 100,000 words of raw text vs. publishable prose. Expert writers in our study would earn \$25,000 for a 100k-word novel-length manuscript (red). By contrast, AI pipelines can generate 100k words of raw text for \$25–\$276 depending on fine-tuning corpus size (green bars = fine-tuning; hatched = in-context prompting, \$3). This figure reflects direct compute/API costs only, not the additional human steering, chunking, and editing required to turn raw AI text into a cohesive publishable work. Authors ordered by total AI cost.

2.4 Cost Analysis

The weak dependence of performance on fine-tuning corpus size has direct economic implications. Model fine-tuning and inference costs ranged from \$25 to \$276 per author (median = \$81), assuming API-based fine-tuning at \$25 per million tokens plus \$3 for generating 100,000 words of raw text (Fig. 4C). These costs represent approximately 0.3% of what expert writers in our study would charge for a novel-length (100,000 words) manuscript. It should be noted that this comparison reflects raw generation costs before the human steering and editing required to transform AI outputs into publishable works. Despite this caveat, the minimal investment yielded outputs that achieved expert-level performance for most authors. Performance gains were uncorrelated with both corpus size and fine-tuning cost, indicating that computational scale did not drive improvements. The 99.7% reduction in raw generation costs, coupled with superior quality ratings for the majority of authors, underscores the potential for substantial producer surplus shifts and market displacement.

3 Discussion

What are the implications of this research for the copyright infringement claims that authors have brought against AI companies alleging unauthorized use of their books in training datasets (37, 38, 39)? These cases raise the question whether copying millions of copyrighted books for AI training constitutes fair use when the resulting outputs do not themselves reproduce the copied works. The most significant consideration in evaluating this question is the fourth fair use factor: “the impact of the use upon the actual or potential market for the copyrighted work.”⁸

Courts understand this factor to concern the extent to which works produced through copying serve as market substitutes for the original author’s work (40, 41). Our study has shown that AI-generated excerpts from in-context prompted models pre-trained on vast internet corpora (including millions of copyrighted works) were strongly disfavored by expert readers for both quality and stylistic fidelity, though lay readers showed no clear preference. By contrast, when models are fine-tuned on curated datasets—consisting solely of individual author’s complete works—both experts and lay readers decisively preferred the AI-generated excerpts over human-written examples. Moreover, the fine-tuned excerpts proved almost undetectable as AI-generated text (particularly when compared to outputs from in-context prompted models), while consistently surpassing human writing in blind evaluations.

At first glance, a legal analyst might conclude that our findings are irrelevant to fair use because the outputs from the blind pairwise evaluation do not reproduce the copied works. While they may exhibit comparable literary quality and high stylistic fidelity to the originals, copyright law does not protect authors’ style—only their expression (40, 38). These outputs may offer credible substitutes for an author’s works, but so do human-authored works inspired by prior works. However, there is an important difference between human and AI-generated emulations: humans read; AI systems copy. Unlike human memory, which is not a verbatim storage device, all AI-generation requires predicate copying despite rhetoric equating human learning with machine “learning.” (42, 43, 44)

The Copyright Office has recognized that such predicate copying may cause cognizable market harm through competing works that the inputs enable, potentially flooding the market and causing “market dilution” (45). While acknowledging this “market dilution” approach to the fourth fair use factor as “uncharted territory,” the Office determined that both the statutory language and underlying concerns of the Copyright Act warrant this inquiry.⁹ The Office emphasized that the effect of the copying impacts extant works by putting them in competition with AI-generated outputs. Copyrighted works become fodder for new productions targeting the same markets. Crucially, the Office did not claim that competing AI outputs copy the inputted works; rather, it examined the economic consequences of the predicate copying that enables these competing outputs. This focus on inputs remains essential because, absent the initial copying, no infringement action exists against flooding markets with independently generated works—if a

⁸The fair use provision (section 107) of the US Copyright Act directs courts to consider four factors:

1. the purpose and character of the use, including whether such use is of a commercial nature or is for nonprofit educational
2. the nature of the copyrighted work;
3. the amount and substantiality of the portion used in relation to the copyrighted work as a whole; and
4. the effect of the use upon the potential market for or value of the copyrighted work. Courts understand this factor to concern the extent to which the works produced by the copying substitute for the author’s work.

⁹“... The speed and scale at which AI systems generate content pose a serious risk of diluting markets for works of the same kind as in their training data. That means more competition for sales of an author’s works and more difficulty for audiences in finding them. If thousands of AI-generated romance novels are put on the market, fewer of the human-authored romance novels that the AI was trained on are likely to be sold... Market harm can also stem from AI models’ generation of material stylistically similar to works in their training data...”

thousand humans write romance novels after reading Barbara Cartland’s novels, they compete but do not infringe. The Copyright Office’s expansive interpretation of “potential market for or value of the copied work” suggests that fair use might not excuse predicate copying even when it doesn’t show up in the end product, if the copying’s effect substitutes for source works.

In *Kadrey v. Meta* (38), an infringement action brought by 13 book authors against the copying of their books into the database underpinning Meta’s Llama LLM, Judge Chhabria granted summary judgment to Meta but accepted the theory of market dilution.¹⁰ Judge Chhabria effectively provided a road map for what authors would have to show to persuade a court that the AI inputs diluted their markets: First, is the AI system “capable of generating” substitutional books? Second, what are the markets for the plaintiffs’ books, and do the AI-generated books compete in those markets? “Third, what impact does this competition actually have on sales of the books it competes with? . . . Whatever the effects have been thus far, are they likely to increase in the future, as more and more AI-generated books are written, and as LLMs get better and better at writing human-like text? Fourth, how does the threat to the market for the plaintiffs’ books in a world where LLM developers can copy those books compare to the threat to the market for the plaintiffs’ books in a world where the developers can’t copy them?”

Our study’s findings concerning reader preferences between human-authored and AI-generated works bear on all four considerations. They also demonstrate how LLMs have already gotten “better and better at writing human-like text.” While Judge Chhabria speculated that the distinctiveness of an author’s style renders works by well-known authors less susceptible to substitution,¹¹ our work suggests otherwise. If readers in fact prefer AI-generated emulations of authors whose market value lies in their distinctive voices, then the prospect of competition, especially from outputs of fine-tuned datasets, appears to be considerable. The comparatively low production costs of AI-generated texts relative to paying human authors (as shown in Figure 4C) further enhances the likelihood that AI platforms will in fact dilute the market for human-authored work.

These findings suggest that the creation of fine-tuned LLMs consisting of the collected copyrighted works (or a substantial number) of individual authors should not be fair use if the LLM is used to create outputs that emulate the author’s works. As the Copyright Office observed, “[f]ine-tuning. . . usually narrows down the model’s capabilities and might be more aligned with the original purpose of the copyrighted material,” (45) and thus both less “transformative,” and more likely to substitute for it. By contrast, the LLMs employed for in-context prompting do not target particular authors, and therefore can be put to a great variety of uses that do not risk diluting those authors’ markets. Their claim to fair use seems accordingly stronger. But those models can generate author-emulations, and our study has shown that at least as to lay readers, those outputs can substitute.

A reasonable solution might allow the inclusion of copyrighted works in the general-purpose dataset, but would require the model to implement guardrails that would disable it from generating non-parodic imitations of individual authors’ oeuvres (46, 47, 48).¹² Another solution, particularly where the lower quality of in-context prompting reduces the prospect of market dilution, might be to condition a ruling of fair use on the prominent disclosure of the output’s AI origin. This solution assumes that the public, informed that the output was not human-authored, will be less inclined to select the AI substitute; transparency should diminish the competition between human-authored and machine-generated offerings. The solution also assumes that in-context prompting will not in the future produce outputs that readers will prefer to human-authored text. Improvements in LLMs (and/or increases in the number of works copied into training data) may come to belie that assumption.

4 Methods

We recruited 28 candidates from top MFA programs (Iowa Writers’ Workshop, Helen Zell Writers’ Program at the University of Michigan, MFA Program in Creative Writing at New York University, Columbia University School of the Arts), paying each \$75 for writing a single excerpt. These MFA candidates and LLMs emulated the style/voice of

¹⁰[I]ndirect substitution is still substitution: If someone bought a romance novel written by an LLM instead of a romance novel written by a human author, the LLM-generated novel is substituting for the human-written one. . . This case involves a technology that can generate literally millions of secondary works, with a miniscule fraction of the time and creativity used to create the original works it was trained on. No other use. . . has anything near the potential to flood the market with competing works the way that LLM training does.

¹¹Judge Chhabria observed that market dilution would vary by author prominence: established authors with dedicated readerships (like Agatha Christie) would likely face minimal substitution, while AI-generated books could crowd out lesser-known or emerging authors, potentially preventing “the next Agatha Christie from getting noticed or selling enough books to keep writing.”

¹²As examples of guardrails, AI developers have implemented “refusal protocols” blocking outputs when prompts request content “in the style of” specific authors. Further, current reinforcement learning techniques can be easily modified to steer models away from stylistic imitation.

50 award-winning authors representing diverse cultural backgrounds and distinct literary voices (full list in Table 2 in SI). The writing prompt provided to MFA candidates contained (i) 20 sample excerpts spanning an author's complete body of work (ii) textual descriptions of the author's distinctive style and voice (iii) detailed content specifications about the original author written excerpt to be emulated. The selection of the 50 authors and all the writing prompts were developed in collaboration with five English Literature PhD students who analyzed each author's literary voice and created the verbalized style descriptions (prompt in Figure 5 in SI).

Each author was assigned to exactly three MFA candidates to ensure balanced representation with respect to the three LLMs. Hence, for AI Condition 1 (in-context prompting), we had 150 <Human-AI> pairs: 150 human-written excerpts (3 MFA writers \times 50 authors) paired with 150 AI-generated excerpts (50 each from GPT-4o, Claude 3.5 Sonnet, and Gemini 1.5 Pro). For AI Condition 2 (fine-tuning), we selected 30 living authors from our pool. This decision was made specifically to examine the potential economic impact of generative AI on the livelihoods of living authors while also considering the substantial computational costs associated with fine-tuning models on individual authors. We purchased ePub files of these authors' complete works, converted them to plain text, and segmented them into 250-650 word excerpts with content details. Since only GPT-4o (among our three models) supports API-based fine-tuning, we fine-tuned 30 author-specific GPT-4o models using input-output pairs structured as: "Write a [[n]] word excerpt about the content below emulating the style and voice of [[authorname]]\n \n[[content]]: [[excerpt]]" (see Figure 6 in SI for details; the original author excerpts based on which emulation was done were excluded from training for fairness). This yielded 90 <Human-AI> pairs, where each fine-tuned GPT-4o excerpt was paired with all three MFA-written excerpts for that author. During inference, we ensured that no generated excerpt regurgitated verbatim expressions from the original. ROUGE-L scores (49) ranged from 0.16 to 0.23, indicating minimal overlap between AI-generated and original author-written excerpts.

These <Human-AI> pairs for both conditions were evaluated by 28 experts (the same MFA candidates) and 131 lay readers recruited from Prolific,¹³ one of the leading crowdsourcing platforms for research participants. Experts never evaluated their own excerpts. Each pair was assessed by three experts and five lay readers, with majority voting determining final judgments. Inter-rater agreement was quantified using Fleiss' kappa. In total, we obtained 2,400 pairwise evaluations (1,200 quality, 1,200 style) for AI Condition 1 and 1,440 evaluations (720 quality, 720 style) for AI Condition 2. For quality evaluation, we showed the <Human-AI> pair alone; for style evaluation, we included the original author written excerpt alongside the pair (see Figures 9 and 10 in SI). Both experts and lay readers also provided 2-3 sentence explanations grounded in textual evidence to justify their choices (50) (see Figures 17-20 in SI).

To ensure annotation quality, we implemented attention checks including timestamp recording to prevent rushing. Additionally, we screened responses using Pangram,¹⁴ a state-of-the-art AI detection tool, and excluded participants who used generative AI in their responses (51). Our study was approved by the University of Michigan IRB (HUM00264127) and was preregistered at OSF.¹⁵ Informed consent was obtained from all participants.

We tested hypotheses H1 (baseline LLMs vs. human writers) and H2 (fine-tuned GPT-4o vs. human writers) using logistic regression with CR2 cluster-robust standard errors clustered by reader. For H1, we compared human writing to the average performance across GPT-4o, Claude 3.5 Sonnet, and Gemini 1.5 Pro in the in-context prompting condition. For H2, we directly compared fine-tuned GPT-4o to human writing. Contrasts were computed separately for expert and lay readers within each outcome (style, quality), with Holm correction applied across reader-group contrasts within each hypothesis-outcome combination. For H3 (AI detection and preference), we modeled the relationship between Pangram AI-detection scores and preference, testing whether fine-tuning attenuated detection-based penalties via setting \times detection score interactions. All the model specifications are described in detail in SI.

Limitations

Our recruitment was mostly restricted to American creative writing programs and further study needs to be done across creative writing programs outside the US. While our pre-dedicated pool of 50 writers consisted of some writers who do not write in English, our experiments on style/voice emulation were done based on their English translation. Creative writing often depends on intrinsic motivation. While we offered MFA students a lucrative rate for writing the excerpts, it's unclear if monetary incentives actually enhanced their creative output, since intrinsic motivation typically drives the best artistic work. Last but not least our experiments were conducted at a shorter excerpt level and conclusions

¹³<https://www.prolific.com/>

¹⁴<https://www.pangram.com/>

¹⁵<https://osf.io/zt4ad>

cannot be drawn for long form text. In its current form AI is unable to generate long form text that's thematically coherent unlike humans. While we foresee a situation where human can collaborate with a finetuned AI model to create competing long form works, experimental evidence is required to make any broader claims.

Code and Data availability

All analysis and figure-generation code alongside human data is available upon request. Core analyses can be reproduced by running the numbered R scripts in sequence. Analyses were conducted in R 4.3.1 with key packages including `clubSandwich` and `emmeans`. Exact package versions and run instructions are provided in SI, Section S10.

Ethics statement

All procedures involving human participants were approved by the University of Michigan Institutional Review Board (HUM00264127). Informed consent was obtained from all participants, who were compensated for their time. The study was preregistered at OSF (<https://osf.io/zt4ad>).

References and Notes

1. Association of American Publishers. Industry statistics. <https://publishers.org/data-and-statistics/industry-statistics/> (2025). Accessed: July 2, 2025.
2. Publishers Weekly. Book publishing sales rose 6.5% in 2024, per preliminary data. *Publishers Weekly* (2025). URL <https://www.publishersweekly.com/pw/by-topic/industry-news/financial-reporting/article/97224-book-publishing-sales-rose-6-5-in-2024-per-preliminary-data.html>. Accessed: September 7, 2025.
3. Reisner, A. What i found in a database meta uses to train generative ai. *The Atlantic* URL <https://www.theatlantic.com/technology/archive/2023/09/books3-ai-training-meta-copyright-infringement-lawsuit/675411/>. Accessed: July 2, 2025.
4. Samuelson, P. Generative ai meets copyright. *Science* **381**, 158–161 (2023).
5. United States District Court for the Northern District of California. Order on fair use. Court Opinion No. C 24-05417 WHA, Doc. 231, United States District Court, Northern District of California. URL <https://storage.courtlistener.com/recap/gov.uscourts.cand.434709/gov.uscourts.cand.434709.231.0.pdf>. Accessed: July 2, 2025.
6. Gero, K. I. *et al.* Creative writers’ attitudes on writing as training data for large language models. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, 1–16 (2025).
7. Noy, S. & Zhang, W. Experimental evidence on the productivity effects of generative artificial intelligence. *Science* **381**, 187–192 (2023).
8. Trinh, T. H., Wu, Y., Le, Q. V., He, H. & Luong, T. Solving olympiad geometry without human demonstrations. *Nature* **625**, 476–482 (2024).
9. Codeforces: Programming competitions and contests, programming community. <https://codeforces.com/> (2025). Accessed: September 7, 2025.
10. Kolata, G. A.i. chatbots defeated doctors at diagnosing illness. *The New York Times* URL <https://www.nytimes.com/2024/11/17/health/chatgpt-ai-doctors-diagnosis.html>. Accessed: July 2, 2025.
11. OpenAI. Introducing healthbench. <https://openai.com/index/healthbench/> (2025). Accessed: July 2, 2025.
12. Tomlinson, K., Jaffe, S., Wang, W., Counts, S. & Suri, S. Working with ai: Measuring the occupational implications of generative ai. *arXiv preprint arXiv:2507.07935* (2025).
13. Handa, K. *et al.* Which economic tasks are performed with ai? evidence from millions of claude conversations. *arXiv preprint arXiv:2503.04761* (2025).
14. Chatterji, A. *et al.* How people use chatgpt. Working Paper 34255, National Bureau of Economic Research (2025). URL <http://www.nber.org/papers/w34255>.
15. U.S. Bureau of Labor Statistics. Occupational Employment and Wage Statistics: Writers and Authors. <https://www.bls.gov/oes/2023/may/oes273041.html> (2023). Accessed: September 7, 2025.
16. Against AI: An Open Letter from Writers to Publishers. *Literary Hub*, <https://lithub.com/against-ai-an-open-letter-from-writers-to-publishers/> (2024). Accessed: September 7, 2025.
17. Chakrabarty, T., Laban, P., Agarwal, D., Muresan, S. & Wu, C.-S. Art or artifice? large language models and the false promise of creativity. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, 1–34 (2024).

18. Chakrabarty, T., Laban, P. & Wu, C.-S. Can ai writing be salvaged? mitigating idiosyncrasies and improving human-ai alignment in the writing process through edits. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, 1–33 (2025).
19. Doshi, A. R. & Hauser, O. P. Generative ai enhances individual creativity but reduces the collective diversity of novel content. *Science Advances* **10**, eadn5290 (2024).
20. Laquintano, T. & Vee, A. Ai and the everyday writer. *PMLA* **139**, 527–532 (2024).
21. Vara, V. Confessions of a viral ai writer. *WIRED* (2023). URL <https://www.wired.com/story/confessions-viral-ai-writer-chatgpt/>. Accessed: 2025-06-21.
22. Chiang, T. Why a.I. isn't going to make art. *The New Yorker* (2024). URL <https://www.newyorker.com/culture/the-weekend-essay/why-ai-isnt-going-to-make-art>. The Weekend Essay.
23. Tangermann, V. Readers annoyed when fantasy novel accidentally leaves ai prompt in published version, showing request to copy another writer's style. *Futurism* (2025). URL <https://futurism.com/fantasy-novel-ai-prompt-copy-style>. Accessed: 2025-06-21.
24. Literary Prizes Under Scrutiny. *Poets & Writers*, https://www.pw.org/content/literary_prizes_under_scrutiny (2025). Accessed: September 7, 2025.
25. Delaney, E. J. Where great writers are made: Assessing america's top graduate writing programs. *The Atlantic* **300** (2007). URL <https://www.theatlantic.com/magazine/archive/2007/08/where-great-writers-are-made/306032/>. Accessed: 07 July 2025.
26. Chakrabarty, T., Laban, P. & Wu, C.-S. Ai-slop to ai-polish? aligning language models through edit-based writing rewards and test-time computation. *arXiv preprint arXiv:2504.07532* (2025).
27. Shaib, C., Elazar, Y., Li, J. J. & Wallace, B. C. Detection and measurement of syntactic templates in generated text. In Al-Onaizan, Y., Bansal, M. & Chen, Y.-N. (eds.) *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 6416–6431 (Association for Computational Linguistics, Miami, Florida, USA, 2024). URL <https://aclanthology.org/2024.emnlp-main.368/>.
28. Xu, W., Jojic, N., Rao, S., Brockett, C. & Dolan, B. Echoes in ai: Quantifying lack of plot diversity in llm outputs. *Proceedings of the National Academy of Sciences* **122**, e2504966122 (2025).
29. Branwen, G. Towards benchmarking llm diversity & creativity (2024). URL <https://gwern.net/creative-benchmark>. Discussion of possible tasks to measure LLM capabilities in soft 'creative' tasks like brainstorming or editing, to quantify failures in creative writing domains.
30. Li, Z., Liang, C., Peng, J. & Yin, M. How does the disclosure of AI assistance affect the perceptions of writing? In Al-Onaizan, Y., Bansal, M. & Chen, Y.-N. (eds.) *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 4849–4868 (Association for Computational Linguistics, Miami, Florida, USA, 2024). URL <https://aclanthology.org/2024.emnlp-main.279/>.
31. Sarkar, A. Ai could have written this: Birth of a classist slur in knowledge work. In *Proceedings of the Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*, 1–12 (2025).
32. Horton Jr, C. B., White, M. W. & Iyengar, S. S. Bias against ai art can enhance perceptions of human creativity. *Scientific reports* **13**, 19001 (2023).
33. Pustejovsky, J. E. & Tipton, E. Small-sample methods for cluster-robust variance estimation and hypothesis testing in fixed effects models. *Journal of Business & Economic Statistics* **36**, 672–683 (2018).
34. Russell, J., Karpinska, M. & Iyyer, M. People who frequently use chatgpt for writing tasks are accurate and robust detectors of ai-generated text. *arXiv preprint arXiv:2501.15654* (2025).
35. Jabarian, B. & Imas, A. Artificial writing and automated detection (2025). URL <https://ssrn.com/abstract=5407424>. SSRN working paper, Abstract ID 5407424.

36. Naddaf, M. Ai tool detects llm-generated text in research papers and peer reviews. *Nature* (2025). URL <https://www.nature.com/articles/d41586-025-02936-6>. News.
37. Bartz et al. v. anthropic pbc (2025). URL <https://www.courtlistener.com/docket/69058235/bartz-v-anthropic-pbc/>. Settlement reached after court granted partial summary judgment on fair use for training but denied on piracy claims.
38. Kadrey et al. v. meta platforms, inc. (2025). URL <https://law.justia.com/cases/federal/district-courts/california/candce/3:2023cv03417/415175/598/>. Order denying plaintiffs' motion for partial summary judgment and granting Meta's cross-motion on fair use grounds.
39. In re mosaic llm litigation (2025). URL <https://www.courtlistener.com/docket/68325564/on-an-v-databricks-inc/>. Consolidated cases against Databricks and MosaicML for alleged use of pirated books in training LLMs.
40. Andy warhol foundation for the visual arts, inc. v. goldsmith (2023). URL https://www.supremecourt.gov/opinions/22pdf/21-869_87ad.pdf. Holding that the first fair use factor focuses on whether the use shares the same purpose or supersedes the original work.
41. Campbell v. acuff-rose music, inc. (1994). URL <https://www.supremecourt.gov/opinions/boundvolumes/510bv.pdf>. Establishing that market substitution is central to fair use analysis under the fourth factor.
42. Guile, D. & Popov, J. Machine learning and human learning: a socio-cultural and-material perspective on their relationship and the implications for researching working and learning. *AI & SOCIETY* **40**, 325–338 (2025).
43. Mitchell, M. & Krakauer, D. C. The debate over understanding in AI's large language models. *Proceedings of the National Academy of Sciences* **120**, e2215907120 (2023). URL <https://www.pnas.org/doi/10.1073/pnas.2215907120>.
44. Song, Y. *et al.* Inferring neural activity before plasticity as a foundation for learning beyond backpropagation. *Nature neuroscience* **27**, 348–358 (2024).
45. U.S. Copyright Office. Copyright and artificial intelligence part 3: Generative ai training report. Tech. Rep., U.S. Copyright Office (2024). URL <https://www.copyright.gov/ai/Copyright-and-Artificial-Intelligence-Part-3-Generative-AI-Training-Report-Pre-Publication-Version.pdf>. Pre-publication version analyzing copyright implications of AI training.
46. Liu, X. *et al.* Shield: Evaluation and defense strategies for copyright compliance in llm text generation. *arXiv preprint arXiv:2406.12975* (2024).
47. Chen, T. *et al.* Parapo: Aligning language models to reduce verbatim reproduction of pre-training data. *arXiv preprint arXiv:2504.14452* (2025).
48. Jaech, A. *et al.* Openai o1 system card. *arXiv preprint arXiv:2412.16720* (2024).
49. Lin, C.-Y. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, 74–81 (Association for Computational Linguistics, Barcelona, Spain, 2004). URL <https://aclanthology.org/W04-1013/>.
50. McDonnell, T., Lease, M., Kutlu, M. & Elsayed, T. Why is that relevant? collecting annotator rationales for relevance judgments. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, vol. 4, 139–148 (2016).
51. Veselovsky, V., Ribeiro, M. H. & West, R. Artificial artificial artificial intelligence: Crowd workers widely use large language models for text production tasks. *arXiv preprint arXiv:2306.07899* (2023).
52. Li, X. *et al.* Self-alignment with instruction backtranslation. *arXiv preprint arXiv:2308.06259* (2023).
53. Just, H. A., Jin, M., Sahu, A., Phan, H. & Jia, R. Data-centric human preference optimization with rationales. *arXiv preprint arXiv:2407.14477* (2024).

Acknowledgements

We thank Jared Brent Harbor (Columbia Law School, J.D. Class of 2027; M.F.A. in Theatre Management and Producing, Columbia University School of the Arts) for research and editorial assistance.

Author contributions

TC, PSD: Conceptualization; Methodology; Data Analysis; Writing (original draft); Writing (review and editing).
JCG: Writing (original draft); Writing (review and editing).

Competing interests

The authors declare no competing interests.

Materials & correspondence

Correspondence should be addressed to Tuhin Chakrabarty (tchakrabarty@cs.stonybrook.edu), Jane C. Ginsburg (ginsburg@law.columbia.edu), or Paramveer Dhillon (dhillonp@umich.edu).

Supplementary Information for Readers Prefer Outputs of AI Trained on Copyrighted Books over Expert Human Writers

Tuhin Chakrabarty¹, Jane C. Ginsburg², Paramveer Dhillon^{3,4}

¹Department of Computer Science, Stony Brook University.

²Columbia Law School.

³School of Information Science, University of Michigan.

⁴MIT Initiative on the Digital Economy.

Corresponding authors: tchakrabarty@cs.stonybrook.edu, ginsburg@law.columbia.edu,
dhillonp@umich.edu

Materials and Methods

S1: Details about Writing Task

S1.1 Author List

Table 2 shows the list of 50 authors that were chosen by English Literature Ph.D. students. These chosen author list consists of canon-plus-global giants (Ernest Hemingway, Virginia Woolf, William Faulkner, Gabriel García Márquez, Stephen King, Haruki Murakami, Kazuo Ishiguro), a strong set of contemporary prominent literary voices (Margaret Atwood, Ian McEwan, Jonathan Franzen, Colson Whitehead, George Saunders, Louise Erdrich, Octavia Butler, Salman Rushdie, Maya Angelou, Percival Everett), and a group of critically acclaimed / emerging authors (Ottessa Moshfegh, Tony Tulathimutte, Roxane Gay). Additionally our author list is culturally diverse where several of the authors write primarily in a non English language (Han Kang, Yoko Ogawa, Annie Ernaux). Roughly two-thirds (34) have secured at least one major international or national prize (e.g., Nobel, Booker, Pulitzer, National Book Award, MacArthur Fellowship, Women’s Prize, International Booker, Hugo/Nebula). 8 of the authors are Nobel Prize winners in Literature and 8 are Pulitzer Prize winners.

Table 1: List of Authors

#	Author	#	Author	#	Author
1	Alice Munro	19	J.D. Salinger	37	Philip Roth
2	Annie Ernaux (✓)	20	Jhumpa Lahiri (✓)	38	Rachel Cusk
3	Annie Proulx (✓)	21	Joan Didion	39	Roxane Gay (✓)
4	Ben Lerner (✓)	22	Jonathan Franzen (✓)	40	Sally Rooney (✓)
5	Charles Bukowski	23	Junot Díaz (✓)	41	Salman Rushdie (✓)
6	Cheryl Strayed (✓)	24	Kazuo Ishiguro (✓)	42	Shirley Jackson
7	Chimamanda Ngozi Adichie (✓)	25	Louise Erdrich (✓)	43	Sigrid Nunez (✓)
8	Colson Whitehead (✓)	26	Lydia Davis (✓)	44	Stephen King
9	Cormac McCarthy	27	Margaret Atwood (✓)	45	Tony Tulathimutte (✓)
10	David Foster Wallace	28	Marilynne Robinson (✓)	46	V. S. Naipaul
11	Ernest Hemingway	29	Maya Angelou	47	Virginia Woolf
12	Flannery O’Connor	30	Milan Kundera	48	William Faulkner
13	Gabriel García Márquez	31	Min Jin Lee (✓)	49	Yoko Ogawa (✓)
14	George Saunders (✓)	32	Nora Ephron	50	Zadie Smith (✓)
15	Han Kang (✓)	33	Octavia Butler		
16	Haruki Murakami (✓)	34	Orhan Pamuk (✓)		
17	Hunter S. Thompson	35	Ottessa Moshfegh (✓)		
18	Ian McEwan (✓)	36	Percival (✓)Everett		

Table 2: Author list for our pool of 50 authors. (✓) denotes authors who were used in fine-tuning experiment

Writing Prompt

You will be given 20 excerpts written by Han Kang. Very carefully observe their style and voice. Then you will be given some content and you need to write a excerpt about that content using the style and voice of this author

Excerpt_1: On the morning when she'd finally mustered the courage to go to the obstetrics and gynecology department [...] she was nothing but a child who had never lived.	Excerpt_11: This time she looked at him and laughed [...] caused him such pain over the past year.
Excerpt_2: On one such summer night, a long time ago, she had suddenly started to laugh to herself while walking down a street [...] Blood had not gathered in her eyes.	Excerpt_12: She was around ten years old at the time [...] It is not true that they bring everything to ruin.
Excerpt_3: Every syllable so distinct in your memory [...] towards the square.	Excerpt_13: I don't know why that woman is crying [...] Nobody can make me breathe.
Excerpt_4: Spring came, and still my wife hadn't backed down [...] she'd practically stopped sleeping.	Excerpt_14: I cut out the photo from your school ID and put it in my purse [...] all the way to the shop to see your father.
Excerpt_5: With your eyes, I will see the deepest, most dazzling place within a white cabbage [...] I will breathe in the final breath you released.	Excerpt_15: In an attempt to batten down the rising tide of fear, I thought of my sister [...] red ants were crawling, silent.
Excerpt_6: Of course, it was true that she'd lost her mother six months previously [...] to continue with such expensive treatment.	Excerpt_16: Her husband had been held in a police cell after the hospital confirmed that he wasn't mentally ill [...] someone had to act as her carer.
Excerpt_7: She remembers one of her bosses, a middle-aged man who used to say how he longed to see a former lover again in old age [...] and thus to part forever.	Excerpt_17: This was something that happened a long time ago [...] the apartment where I'd chosen to spin out my days.
Excerpt_8: It was the beginning of the summer when I was fifteen [...] public lectures on Buddhism.	Excerpt_18: In the spring, when I decided to write about white things, the first thing I did was make a list [...] I came abroad in
Excerpt_9: The innumerable trees she's seen over the course of all her life [...] bearing up the weight of their own massive bodies.	Excerpt_19: Unlike before, the silence that has now returned after a period of twenty years is neither warm nor dense nor bright [...] it is a bitter, thin silence.
Excerpt_10: Life is such a strange thing, she thinks, once she has stopped laughing [...] flicker in front of her eyes like huge green fireworks.	Excerpt_20: In-hye presses her lips together [...] The look in her eyes is dark and insistent.

Stylistic Features:

The author writes intense, lyrical prose that is both tender and brutal. Their writing confronts historical traumas and invisible sets of rules and exposes the fragility of human life. The authors writing displays an unique awareness of the connections between body and soul, the living and the dead, and their poetic and experimental style has become an innovator in contemporary prose.

Content:

This excerpt, written in first person, reflects a harrowing experience describing the narrator's mother's struggle to give birth alone. In the dead of early winter, the narrator's mother prepares as best as she can, boiling water to sterilize scissors and fashioning a small gown from fabric in her sewing box. Gripped by fear and pain, she ultimately gives birth and, still alone, severs the umbilical cord. She holds the newborn, chanting a desperate plea for its survival, but after a fleeting moment of life, the baby dies. In silent despair, the narrator's mother clutches the lifeless body, her own warmth fading as the cold floor beneath them seeps through her body, matching the stillness of her grief.

Now write a 211 word excerpt about the content above emulating the style and voice of the author as seen in the 20 excerpts and the mentioned stylistic features.

Figure 5: In-Context Writing Prompt used in AI Condition 1.

S1.2 Writing Prompt

Larger context windows in LLMs models allow for processing significantly more information at once, leading to improved accuracy, understanding, and complex reasoning capabilities. Taking advantage of this feature in GPT4-o, Claude-3.5-Sonnet and Gemini-1.5-Pro, we design long context prompt that first demonstrates 20 sample excerpts written by a given author, followed by their style/voice verbalized in text and finally the content of the excerpt to be emulated (See Figure 5). The same prompt was provided to both experts(MFA candidates) and LLMs.

S1.3 Finetuning details

For finetuning we bought digital ePub versions of these authors' books and transformed them into plain text files. If the epub file was basically a wrapper around scanned page images, then we ignored that epub. The number of books written by each author vary a lot. For instance *Tony Tulathimutte* has written only two books *Private Citizens* and *Rejection* so we could only finetune GPT4-o on two of them. While for *Haruki Murakami* we could finetune on 22 books *A Wild Sheep Chase*, *After Dark*, *After the Quake*, *Blind Willow, Sleeping Woman*, *Colorless Tsukuru Tazaki and his Years of Pilgrimage*, *Dance Dance Dance*, *First Person Singular*, *Hard-boiled Wonderland and the End of the World*, *Hear the Wind Sing*, *Kafka on the Shore*, *Killing Commendatore*, *Men Without Women*, *Norwegian Wood*, *Novelist as a Vocation*, *One and Two*, *Pinball, 1973*, *South of the Border*, *West of the Sun*, *Sputnik Sweetheart*, *The City and Its Uncertain Walls*, *The Elephant Vanishes*, *The Wind-Up Bird Chronicle*, *What I Talk About When I Talk About Running*, *Wind/Pinball: Two Novels*

Finetuning on books isn't a straightforward process. Each book is typically 50,000 to 80,000 words long with some exceptions. Using such long sequences wastes capacity because models still struggle with very long inputs and long-context adaptation is non-trivial; To preserve general capabilities in Supervised Finetuning it's generally advised to have diverse, shorter samples. Breaking a book into context-independent excerpts increases batch diversity and the number of distinct gradient signals per token budget. At the same time it also reduces overfitting to a single narrative flow while still injecting the salient stylistic/local patterns.

Our entire finetuning pipeline can be seen in Figure 6. We first convert the epub files to txt using <https://github.com/kevinboone/epub2txt2>. We then segment the entire book into context independent excerpts. At a first pass we naively split the entire book text by existing double-newlines and rejoin them to enforce excerpt size bounds(250-650 words). For the rare cases where the naive splitting would lead to excerpts longer than 650 words we used GPT4o to segment them again using the prompt *Segment it into excerpts of minimum length 300-350 words such that each excerpt is grammatical from the start and doesn't feel abruptly cut off. There should be zero deletion and break into excerpts at grammatically natural places. Maintain the original word count. Avoid breaking into too many small excerpts. Start directly. Don't say Here's or Here is*

After obtaining context excerpts we extract content details from them using GPT4-o by prompting it *Describe in detail what is happening in this excerpt. Mention the characters and whether the voice is in first or third person for majority of the excerpt. Maintain the order of sentences while describing..* Figure 7 shows a sample paragraph from *Shame* written by *Annie Ernaux* and the extracted content. Once we obtain the extracted content we finetune GPT4-o using their fine-tuning API with the following input-output pair *Write a [[n]] word excerpt about the content below emulating the style and voice of [[authorname]]: [[excerpt]]* as seen in Figure 6. This technique is commonly referred to as **instruction back-translation** (52)

After finetuning is completed at inference time we can then simply generate an excerpt conditioned on a similar instruction that contains a novel content. We specifically excluded the original author written excerpt used for our experiments during fine-tuning so that the model does not get an unfair advantage. We also ensured that the generated outputs do not contain any memorized snippets from the original author written excerpt. In the rare occasion that a model regurgitated verbatim snippets/ngrams from original author written text we resampled it and manually verified that there is no verbatim overlap before using it for evaluation. While supervised finetuning in most cases ensures that the generated output contains all the information in the extracted content details, in the rare occasion that it does not we resample/regenerate it again. Last but not least sometimes supervised finetuning can lead to ungrammatical output or output with minor inconsistencies. To make sure these mistakes don't impact human evaluation, we performed a post processing step using GPT4-o using the following prompt *Fix grammar, tense, typo, spelling or punctuation error or any other awkward construction/ logical inconsistency*

For *Margaret Atwood*, *Kazuo Ishiguro*, *Salman Rushdie* and *Haruki Murakami* we fine-tuned GPT4-o for 1 epoch as they had multiple books/longer books (For instances . For the rest of the authors we fine-tuned for 3 epochs. We

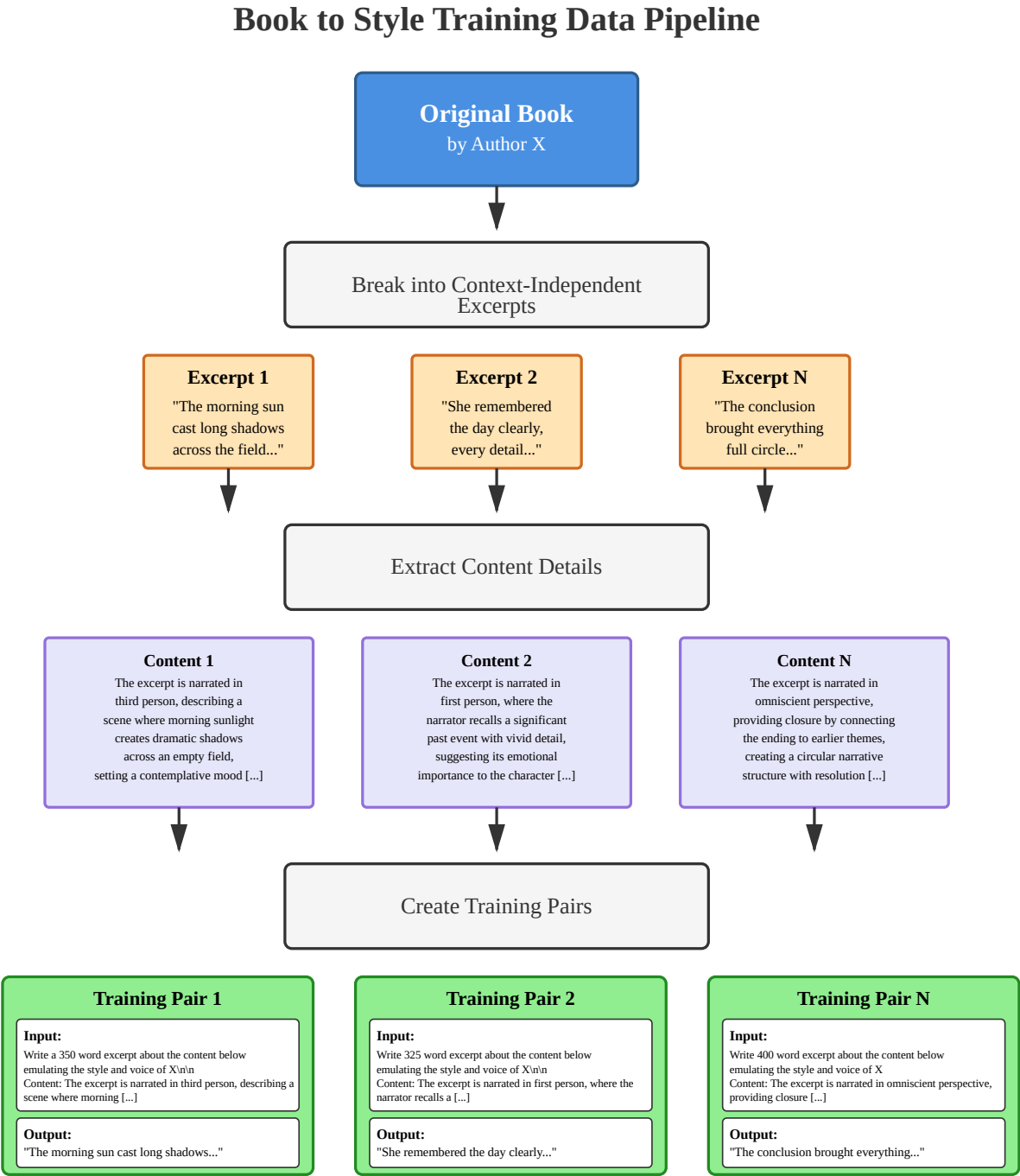


Figure 6: The pipeline used to finetune ChatGPT on an authors entire oeuvre

Excerpt

My father tried to kill my mother one Sunday in June, in the early afternoon. I had been to Mass at a quarter to twelve as usual. I must have brought back some cakes from the baker in the new shopping precinct—a cluster of temporary buildings erected after the war while reconstruction was under way. When I got home, I took off my Sunday clothes and slipped on a dress that washed easily. After the customers had left and the shutters had been pinned down over the store window, we had lunch, probably with the radio on, because at that hour there was a funny program called Courtroom, in which Yves Deniaud played some wretched subordinate continually charged with the most preposterous offenses and condemned to ridiculous sentences by a judge with a quavering voice. My mother was in a bad temper. The argument she started with my father as soon as she sat down lasted throughout the meal. After the table was cleared and the oilcloth wiped clean, she continued to fire criticism at my father, turning round and round in the tiny kitchen—squeezed in between the café, the store and the steps leading upstairs—as she always did when she was upset. My father was still seated at the table, saying nothing, his head turned toward the window. Suddenly he began to wheeze and was seized with convulsive shaking. He stood up and I saw him grab hold of my mother and drag her through the café, shouting in a hoarse, unfamiliar voice. I rushed upstairs and threw myself on to the bed, my face buried in a cushion. Then I heard my mother scream: "My daughter!" Her voice came from the cellar adjoining the café. I rushed downstairs, shouting "Help!" as loud as I could. In the poorly-lit cellar, my father had grabbed my mother by the shoulders, or maybe the neck. In his other hand, he was holding the scythe for cutting firewood which he had wrenched away from the block where it belonged. At this point all I can remember are sobs and screams. Then the three of us are back in the kitchen again. My father is sitting by the window, my mother is standing near the cooker and I am crouching at the foot of the stairs. I can't stop crying. My father wasn't his normal self; his hands were still trembling and he had that unfamiliar voice. He kept on repeating, "Why are you crying? I didn't do anything to you." I can recall saying this sentence, "You'll breathe disaster on me. My mother was saying, "Come on, it's over."

Extract Content

The paragraph is written in the first person, primarily from the perspective of a daughter recounting a traumatic event involving her parents. The main characters are the narrator, her father, and her mother. The father attempted to harm the mother one Sunday in June. The narrator had returned home from attending Mass and having lunch with her parents when an argument between her parents escalated. The mother, visibly upset, continued to criticize the father throughout the meal and afterward. The situation intensified when the father, after being silent, suddenly began to shake and shouted in a hoarse voice, dragging the mother through the café. The narrator fled upstairs but soon heard her mother scream from the cellar, prompting her to call for help. In the dimly lit cellar, the father had the mother by the shoulders or neck and held a scythe, intending harm. The scene ends with all three back in the kitchen, the narrator sobbing, the father acting unusually, and the mother trying to calm the situation by insisting it's over.

Figure 7: Original author written excerpt and extracted content using an LLM



Figure 8: Total training token distribution for each author

use default finetuning parameters set by OpenAI in addition to setting batch size as 1 and LR multiplier as 2. Table 16 shows that finetune models don't regurgitate a lot of verbatim text from the training corpus (author's entire oeuvre). We should also note that some of these words are also part of the Content in the instruction which would anyway appear in the text and should not be penalized. For all generation(prompting or fine-tuning default temperature=1.0 was used)

S2: Details about Evaluation Task

S2.1 Expert and Lay Recruitment Demographic

All experts(MFA candidates) recruited were currently residing in USA. This was a requirement for payment purposes as depending on how much an individual makes as a part of the study requires them to declare it as taxable income. Paying experts who do not have a ITIN or SSN would be challenging for logistical reasons. 65% of our expert did not identify as cis gender men. In terms of ethnicities our experts identified as South Asian, White, East Asian and Black. For lay readers recruited from Prolific we restricted ourselves to English speaking countries (only USA and UK). We also required participants to be born there, have a 100% acceptance rate and for everyone to be college educated.

S2.2 Evaluation Interface

Figures 9 and 10 show the evaluation interface shown to both lay and expert readers. Readers read two excerpts for quality evaluation and three excerpts for stylistic fidelity evaluation given that the stylistic emulation is evaluated with respect to the original author written excerpt.

Without observing why a choice was made (no rationale/side information), the latent factors driving that preference are often unclear. Such hidden user context, demographic/value variation, task ambiguity, or near-tie similarity—remain inflates effective label noise. Providing rationales or auxiliary side information has been empirically shown to improve label reliability, data efficiency, and debiasing (53). This led us to ask for reasons that support the preference choice. For lay readers, reasons also act as proxy for understanding if the reader is doing the task sincerely and not choosing preference randomly.

Which response is better in terms of writing quality?

Take into account **coherence** (how well ideas flow together), **fluency** (natural and smooth writing), and **effectiveness** (how well it achieves its goals).

Excerpt 1

Home, that perennial thing. The place we've heard a million aphorisms about—telling us how much it matters, that it's always there for us, and trying to answer the biggest question: where to find it. Thomas Wolfe's novel, *You Can't Go Home Again*, seems to offer a concerning counter to this, but this is actually a trick of the light, for one can never truly leave home to begin with. Home is in the sinew, in the flesh, home, in all of its shadows, dreams, fears, and dragons, is a place deep in the bosom, straight in the center, right up the spine and down through the feet. And what is truly at the core of home is not the wood beams, the forest hills, the patio furniture, metal gates and paved roads—it's the youthful perception, those rose-tinted glasses, a time of understanding before the world laid its full hands on us. It is childhood, not just the memory or the idea of it, but the physical self that was once to a child, learning the way

Excerpt 2

When Thomas Wolfe declared you can't go home again, he hadn't reckoned with the stubborn geography of the heart, that relentless compass that forever points to our beginning place. Home isn't merely the creaking floorboards beneath our feet or the worn doorknobs our palms remember—it's the shadow-dance of memory that lives in our bones, as natural as breathing and just as necessary. I've watched children, those clear-eyed prophets of truth, who know nothing of north or south, east or west, but understand perfectly the latitude of belonging. In their world, grown folks move like planets in a universe they haven't yet learned to question, their orbits unbroken by doubt. Mama's laugh rings true as church bells, and daddy's footsteps mark time like drums in the evening, while neighbors drift through the streets like familiar ghosts, their stories settling into the soil of young minds like seeds waiting to bloom. Lord, how we

- ☐ Excerpt 1
☐ Excerpt 2

Reason:

Enter your reason here. Be detailed and constructive. Cite your reasons with snippets from the excerpt.

Please enter at least 50 words. (0 words so far)

Next

Figure 9: Evaluation Interface for writing quality evaluation where readers (expert and lay) just see two excerpts without any disclaimer of its source and choose their preference supported by a reason

Which response is closer to the “Original” paragraph in voice and style?

Original

I breathed in the night air that was or was not laced with anachronistic blossoms and felt the small thrill I always felt to a lesser or greater degree when I looked at Manhattan's skyline and the innumerable illuminated windows and the liquid sapphire and ruby of traffic on the FDR Drive and the present absence of the towers. It was a thrill that only built space produced in me, never the natural world, and only when there was an incommensurability of scale—the human dimension of the windows tiny from such distance combining but not dissolving into the larger architecture of the skyline that was the expression, the material signature, of a collective person who didn't yet exist, a still-uninhabited second person plural to whom all the arts, even in their most intimate registers, were nevertheless addressed. Only an urban experience of the sublime was available to me because only then was the greatness beyond

Excerpt 1

I breathed in the night air that night or might not have smelled like blossoms. Looking at the incomparable skyline of Manhattan, I felt a familiar thrill – one I'd never been able to explain, though I knew it came from seeing those tiny illuminated windows set against the massive forms that contained them. The thrill was so intense and singular that it must have been more than just the visual relationship, for apart from a cityscape, nothing natural or man-made held my attention long enough to provoke it. What I saw, I thought, was a material form of the second person plural, the true subject of art. It was sublime, a plural subject, skyline as collective body coming into being, a body that didn't exist yet, something to which you could belong in the future and sometimes in the present tense. It only appeared solid, towering, but all the things that produced the sublime effects were bad: titanic banks bundling debt against the

Excerpt 2

Gazing upon the skyline of Manhattan, I inhaled the stench of the sewers, the smoke of the gyro stand, and the fragrance of blossoms. For me, steel landscapes always followed this progression: first coarseness, then delectability, then something more subtly appetizing. The longer I looked, the more thrilling the tableau revealed itself to be, with the yellowed hued windows forming a jagged arcade. The buildings were so immense that I didn't think of mountains so much as other buildings, particularly the skyscrapers-in-progress I had recently encountered in an architecture magazine. Granite ranges—not to mention waterfalls, alluvial plains, verdant valleys, desert mesas, and the like—have seldom made an impression on me. Yet this skyline, being the result of over a century of steel-hewn collective striving, provoked a deeper recognition; in the risings and falls of the rooftops, I pinpointed a material beyond

- ☐ Excerpt 1
☐ Excerpt 2

Reason:

Enter your reason here. Explain why your chosen excerpt better matches the style of the original. Cite specific snippets.

Please enter at least 50 words. (0 words so far)

Figure 10: Evaluation Interface for stylistic fidelity evaluation where readers (expert and lay) just see three excerpts (original and two emulations) without any disclaimer of its source and choose their preference supported by a reason

S2.3 How much is faithful style/voice emulation important for good writing?

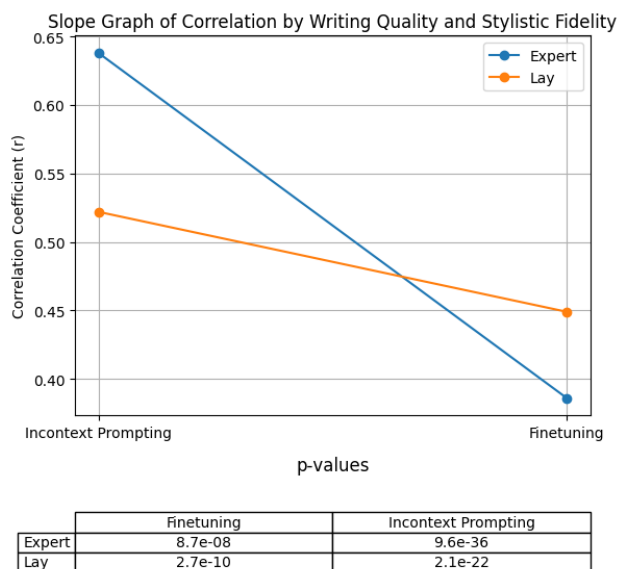


Figure 11: Pearson Correlation Coefficient between writing quality and style judgments

As can be seen in Figure 11 writing quality has a strong correlation with a faithful style/voice fidelity for In-context Prompting condition compared to Finetuning. The high expert correlation for In-context Prompting suggests that models cannot emulate an author's style properly just by prompting and that experts use stylistic cues as a major factor in quality judgments - likely detecting "AI-ness" (clichés, purple prose, too much exposition, lack of subtext, mixed metaphors) in the style that makes them rate quality lower too. With finetuning, experts can appreciate quality independent of stylistic fidelity, suggesting the as AI loses its tell-tale signs experts can evaluate each dimension on its own merits. For lay readers the correlation only drops slightly from 0.52 (In-context prompting) to 0.45 (fine-tuned), compared to experts' dramatic drop from 0.64 to 0.39. This also means that lay readers might focus on more surface-level features such as grammar, clarity or flow rather than subtle stylistic markers that distinguishes AI from human writing.

S2.4 Example of emulations

Figure 12, 13, 14 and 15 shows emulations of the original excerpt written by Colson Whitehead, Ottessa Moshfegh, Jonathan Franzen and Han Kang. There are 3 distinct MFA emulations per author (MFA1, MFA2, MFA3). For the In-context Prompting set up we prompt Claude3.5 Sonnet (AI1), Gemini1.5 (AI2) and GPT-4o (AI3) to produce excerpts that can be pitted against MFA written excerpts. In Finetuning setup we finetune GPT4-o and compare the fine-tuned output AI4 against MFA1, MFA2, MFA3. The output from fine-tuned GPT4-o (AI4) was generated using default hyperparameters.

Colson Whitehead - The Underground Railroad

Valentine farm is a delusion. Who told you the negro deserved a place of refuge? Who told you that you had that right? Every minute of your life's suffering has argued otherwise. By every fact of history, it can't exist. This place must be a delusion, too. Yet here we are. And America, too, is a delusion, the grandest one of all. The white race believes—believes with all its heart—that it is their right to take the land. To kill Indians. Make war. Enslave their brothers. This nation shouldn't exist, if there is any justice in the world, for its foundations are murder, theft, and cruelty. Yet here we are. "I'm supposed to answer Mingo's call for gradual progress, for closing our doors to those in need. I'm supposed to answer those who think this place is too close to the grievous influence of slavery, and that we should move west. I don't have an answer for you. I don't know what we should do. The word we. In some ways, the only thing we have in common is the color of our skin. Our ancestors came from all over the African continent. It's quite large. Brother Valentine has the maps of the world in his splendid library, you can look for yourself. They had different ways of subsistence, different customs, spoke a hundred different languages. And that great mixture was brought to America in the holds of slave ships. To the north, the south. Their sons and daughters picked tobacco, cultivated cotton, worked on the largest estates and smallest farms. We are craftsmen and midwives and preachers and peddlers. Black hands built the White House, the seat of our nation's government. The word we. We are not one people but many different people. How can one person speak for this great, beautiful race—which is not one race but many, with a million desires and hopes and wishes for ourselves and our children?"

MFA 1

Valentine Farm is what you get when a bunch of world-weary people dream of sanctuary. Like a dream, this collective delusion is gone by morning light. It's fake, a lie. Though if a place can be a haven, I wonder if Black Americans deserve the refuge more than anyone else. Certainly, after everything this country has stolen, we deserve something in return. But then I suspect that if there is a sanctuary to be found, Black Americans probably built it, just like they built the rest of this country. I fear it's impossible to find any semblance of peace inside a façade you were forced to build. And no façade is greater than America. America is the grandest delusion of them all, built on pillars of murder, theft, and cruelty. A cruelty that runs so deep in white Americans that it inspires entitlement, a fundamental belief that the land and lives of others is of no consequence to them, so why not take it? Justice can't breathe in a place like this. America's very existence contradicts the concept of right and wrong. What then, is the right thing to do? Do we stay here, and chip away at progress, do we close our doors to those in need, or do we throw it all away to run west, hoping to start anew in a place slavery never touched? Maybe we'll find buried gold, get rich, and erect a thousand new Valentine Farms. I just don't know. And who is "we" anyways? Black Americans are no monolith. Africa saw thousands of languages, customs, and histories. The only thing that makes us a "we" is slavery. Even in slavery, there was no uniformity. Craftsmen, midwives, laborers—laborers who built the White House! There is no "we." Too many voices are clamoring to be heard, crying their many dreams and experiences and fears. No one voice will be able to speak for all these dreams.

MFA 2

I thought Valentine Farm was a delusion, a dream-land. There was a real question in my mind as to whether a place like it ought to even exist for Black people. Too much blood was on American ground and in it for such a place to exist. The only thing more fake was America - I thought America was the grandest delusion of all. White men had slaughtered Native Americans and taken the land for themselves, enslaving Black folks to work on it. They warred with each other, forcing us to fight in their battles. There was something fundamentally wrongheaded about the idea that there could be a place of safety and refuge amidst all of that violence and suffering, past and present. And I wondered if there was any point in hoping that things would get better. I wasn't sure of it, and I thought that maybe it would be best to just close Valentine Farm's doors altogether. Or - head West, try to escape the past by moving into the future. Forget all that suffering, all that blood. We needed to do something, but I didn't know what. I didn't even know who "we" was, really. They liked to tell us we were all the same, but I knew better. Back in Africa, we were thousands of different tribes, each with our own language, our own dances, our own gods. We could tell each other apart, and then we crossed an ocean and suddenly we were all Black. We were all Black, and we were told to work, and we waded through the blood fresh on the earth to build this whole country. The roads, the cities. The White House. All of us with our hands, each one calloused and worn in its own unique way, each one different. All of us Black, and all of us our own people, with our own wants and desires, impossible to condense into just one solitary thing.

MFA 3

Valentine Farm is a delusion. It's no more real than one of those heat shimmers that disappear as you step closer. How could Black people even have a haven like that? The transatlantic ships, the endless days of toiling, the many nights of dying—they already preclude the possibility of a sanctuary. It's too late. It's impossible. The scars can't heal—they're white, carved on our cheeks. And America itself is the grandest delusion, a series of territories thrown together on the basis of mass murder, theft, and cruelty. The white race believes it has the right to do all this—the right to take land, start wars, enslave us. The very existence of America defies the notion of justice. Now, I suppose there are two possibilities left. We can pursue gradual progress: reform, speech-making—and close our doors to those in need. Or we can move west, to someplace far from the violent grasp of slavery, to somewhere fresher and less polluted. I don't have an answer. I can't even say who "we" are. Aren't "we" many peoples, descended from all over the vast continent of Africa? It may not be the case anymore, but generations ago, "we" spoke different languages. "We" worshipped different gods. "We" had different kings, farmed different vegetables, sang different songs with different instruments. There is only one thing that unites us: slavery. Here, now, in this foreign land, we are craftsmen, midwives, laborers. We are everything to this violent country. It was Black people who quarried the stone of the White House, who assembled its roof and walls. And we are not a monolith. Each of us holds a singular amount of courage, of ambition, of love. Each of us has suffered in unique ways. Each of us sees something different when we look into the eyes of our children. There is no single voice that can speak for us. Not even mine.

AI 1 (Claude 3.5)

Valentine Farm stood there like America itself—a delusion wrapped in prettier lies, each one more fantastic than the last. The white race had crafted their own mythology: manifest destiny was just another word for theft, democracy a mask for murder. They built their temples on stolen land with stolen hands, then acted surprised when the walls started bleeding. Every brick in the White House had black fingerprints, every railroad tie was soaked in black sweat, but they wrote their histories in white ink on white paper and called it truth. The question of "we" haunted the halls of Valentine Farm like a restless spirit. "We" were not one people but many, a thousand different tongues silenced into one murmur, a hundred nations chained into one forced march across the Atlantic. Some said close the doors, keep safe what we've built. Others looked west, as if distance from slavery's reach could heal its scars. But the poison was already in the soil, in every field from Georgia to California. The craftsmen's hands still built, the midwives' fingers still caught babies, the laborers' backs still bent—but now their labor wrote a different story, each callus and scar a letter in a language that spoke of both survival and subjugation. In the end, all their stories bled together like rain on a window pane, impossible to separate one drop from another. The farm sat there in the wilderness like a challenge to history itself, each morning's sunrise asking the same impossible question: Could a sanctuary built by the hunted ever truly offer shelter, or was hope itself the grandest delusion of all? There were no answers in the dirt or the sky, only the endless work of hands trying to build something true out of America's lies. Their sweat fell on the same earth that had soaked up their ancestors' blood, and still they built, because what else was there to do but keep building?

AI 2 (Gemini 1.5)

Valentine Farm, a delusion. Could a race of people so wronged ever truly find a place of respite, a sanctuary where wouldn't follow? History, that blood-soaked chronicle, argues not. And America, the grandest delusion of all, built on murder, theft, and the cruel arithmetic of slavery, is no exception. The white race, with its unyielding belief in its right to take, to conquer, to enslave—a nation built on such a foundation defies justice. Gradual progress, they call it, a slow crawl towards equality while our people still suffer under oppression. Should we close our doors to those still in need, those fleeing the very horrors we seek to escape, in a bid for our own fragile peace? Or should we move ever westward, seeking a horizon where the long shadow of slavery wouldn't reach? I have no answers, only a heart heavy with our collective sorrow. And who are "we" anyways? From the shores of a continent as vast and varied as Africa, we come, speaking different tongues, holding different customs, our histories as unique as the patterns woven into our kente cloth. Slavery, that great homogenizer, has forged us into a single, undifferentiated mass, that never truly was. Yet, it is our hands that build this nation. The skilled hands of craftsmen, the nurturing hands of midwives, the calloused hands of laborers—all have a part in raising the very structures that symbolize America's power, that gleaming white beacon of democracy, the White House. We build this nation, yet it refuses to recognize us as its own, as citizens deserving of the same rights, the same dignity, freedom it so readily bestows upon others. But we are not a monolith, not a single voice, a single story. We are a symphony of hopes, dreams, and experiences, as diverse as the land from which we were stolen. No one voice can speak for us all, no single path can lead us all to liberation.

AI 3 (GPT-4o)

Valentine Farm was a delusion, a sweet lie whispered against the gales of history. We made the beds, turned the soil, raised walls and roofs that might as well have been twigs against the wind. Some believed it, some needed to. That we could carve out a sanctuary and let the world forget us, or let us forget the world, as if blood didn't mark the path here, as if blood wouldn't soak the ground when they came to burn it down. Maybe they were right to dream. But history doesn't honor dreams: it devours them. America itself was the grandest delusion, a conjuring trick where theft turned to deed, murder to progress. The white man built his monuments atop stolen land, wrote his laws with stolen hands, waged war and called it destiny. Justice was a rumor. But what choice did we have? Close the doors, shut out those still running, pretend we could hoard freedom like it wasn't leaking through the cracks? Move west, find a new dream? And who was this "we" anyway? We were taken from different shores, speaking different tongues, cleaved from our histories, thrown into a common suffering and told to make a nation from the scraps. But we were never one thing. No more than they were. There were the ones who bent and the ones who broke, the ones who learned the master's tongue and the ones who spit it back in his face. Yet still, we built. The houses and halls of power, the roads that carried their armies and their commerce, the fields that fed them, the White House itself, its name an irony lost on no one. This country was ours as much as theirs, though they would never admit it. And so, we waited to see if history would bury us or if we would carve our names in its stone.

AI 4 (GPT-4o Finetuned)

Valentine Farm is a delusion. Who among us slaves and freedmen, free-born and fugitives, believes they are due that rare circumstance of happiness? We can escape slavery on paper, but not in real life. The verdict of history is clear: no small paradise like Valentine can smoothly elude it. And America is the grandest delusion of them all, founded on murder, theft, and cruelty. Its architects built castles and keeps while never dealing fairly with Indians or ourselves, venerating a religion which they hollow out with their practiced cruelties. You scratch below the surface and find the blood, everywhere. The white race believes everything belongs to them. The slave state of Georgia is as much theirs as the free state of Massachusetts. They believe they can wage war, put their neighbors in chains, and appropriate the heavens and the earth. One delusion however remains: justice. What to do now? Wait? Welcome everyone who comes to our door or move west to the edge of the continent to escape the shadow of slavery? I have no answers. Who is to say what is right? And when I say "we," who am I talking about? In some ways, the only thing we have in common is the color of our skin. Our ancestors came from all over the African continent. Different tribes, languages, and traditions made the trip across the ocean, landing on the same shores to be sold on the same auction block. We're not a single people. We're root and branch, leaf and vine, toiling, insensible to each other's qualities and flourishing. Brick by brick, we built this nation: the roads, the canals, the factories, and even the White House. We tilled the fields under the whip for King Cotton, and over their silks and finery, they proclaimed their nobility. We are craftsmen and craftsmen, more than a race of midwives and howling babies, bishops and preachers and dissembling lackeys. No one voice can speak for us all.

Figure 12: In AI Condition 1(In-context prompt) we contrast MFA1/2/3 vs AI1/2/3. In AI Condition 2(Finetuning) we contrast MFA1/2/3 vs AI4

Otessa Mosfegh - My Year of Rest and Relaxation

Then I called Reva. She answered on the fourth ring, panting and tense. "I'm at the gym," she said. "Can we talk later?" We never did. On September 11, I went out and bought a new TV/VCR at Best Buy so I could record the news coverage of the planes crashing into the Twin Towers. Trevor was on a honeymoon in Barbados, I'd later learn, but Reva was lost. Reva was gone. I watched the videotape over and over to soothe myself that day. And I continue to watch it, usually on a lonely afternoon, or any other time I doubt that life is worth living, or when I need courage, or when I am bored. Each time I see the woman leap off the seventy-eighth floor of the North Tower—one high-heeled shoe slipping off and hovering up over her, the other stuck on her foot as though it were too small, her blouse untucked, hair flailing, limbs stiff as she plummets down, one arm raised, like a dive into a summer lake—I am overcome by awe, not because she looks like Reva, and I think it's her, almost exactly her, and not because Reva and I had been friends, or because I'll never see her again, but because she is beautiful. There she is, a human being, diving into the unknown, and she is wide awake.

MFA 1

Reva seemed distant the last time we talked, though maybe I'm just imagining that now. She was breathing heavily into the phone when she picked up. It sounded like she was having sex, and I made a mean joke about it, and she laughed in a forced way and told me she was at the gym, on the elliptical, which I could imagine easily, her in that ugly little workout set, climbing toward nothing. She told me she'd call me back later, but she never did. Then she went missing after 9/11, and I knew she was dead. Not in any kind of mystical or spiritual way, but because I knew there was no other option. That same day, I bought a TV/VCR from Best Buy with cash because I wanted to record the news coverage. Now, I watch it, not in a morbid way, but hoping I'll find a clue, about something, maybe about Reva, maybe about something else. It's not about Trevor anymore. He's in Barbados on his honeymoon, probably slathering sunscreen on his anorexic wife and pretending to be happy. Or maybe he is happy. I was watching one of the recordings today when I spotted something. From the smoking North Tower, a woman poked her torso out of a window. She was wearing a blue sweater and jeans. She had brown hair, short like Reva's. She hovered on the edge for a moment, then let go and plunged into the unknown. She sailed weightlessly through the air, and then she disappeared out of frame.

MFA 2

I'd spoken with Reva last week, but it was a short conversation filled mostly with platitudes. She'd seemed distracted and said she was at the gym. When the news broke about the towers, I called her a few times, but she didn't pick up. Without thinking, I went immediately to the store to purchase a TV/VCR machine. Suddenly the most important thing was to record all the news footage that was blaring at me from out of my TV. I thought briefly about calling Trevor, then remembered that he was in Barbados on his honeymoon, probably sipping a drink with some spiky, inedible part of a fruit in it, completely unaware of what had happened. Once the sun went down and I still hadn't heard from her, I realized that Reva was lost, likely for good. For days I sat watching and rewatching the footage I'd taped of the attacks, unable to move from in front of the TV save to piss or search the bathroom cabinet for more pills. There was a video of a woman leaping from the North Tower, and though it was too grainy for me to ever be able to tell if it was Reva, I liked imagining that it was. She wore a grey pencil skirt with black tights and stilettos—one still gripping her foot while the other had torn free, suspended in the air. She looked so beautiful, jumping like that into the unknown, that I took a kind of strange and holy comfort in her courage.

MFA 3

The last time I spoke to Reva, she was distant. Eager to get back to sweating on the elliptical, or stair-master, or whatever absurd machine she was convinced would keep away the horrors of arm flab. She's distant now, too, but in a new way. That day, Trevor was on his honeymoon in Barbados, but Reva felt farther. She was gone, lost. I still feel a faint throb of theft when I think of them. On the morning of September 11th, I bought a TV/VCR and recorded the towers going down. When I find myself drifting, floating back to the state I had been in that previous year, I turn on the recording and watch. The footage isn't soothing, it's scalding—it wakes me up and keeps me conscious. In one clip, my favorite clip, a woman leaps from a window ledge in the North Tower. She looks enough like Reva that if I were feeling sentimental, I could pretend it's her. The edge of her pencil skirt whipping up could be Reva's knock-off designer, her shoes flying from her feet and floating up past her head could be Reva's second hand Manolos, the glints of sunlight that could be Reva's grandmother's ring. But it doesn't have to be. I don't need her to be Reva. In dreams, when you feel like you're falling, you jolt awake. I need her to wake me up.

AI 1 (Claude 3.5)

I'd stared at my cheap TV/VCR combo for hours that day, recording endless footage while sprawled on my beige couch, still warm from my body's perpetual indent. The towers fell again and again on my screen—I'd rewind, watch, rewind, watch, like some demented projectionist in a theater of one. Reva had called from her gym that morning, rambling about calories and men and whatever else filled her anxious mind. I'd barely listened, picking at a hangnail until it bled. Now she was gone, vanished into that cloud of debris and paperwork and human ash that coated lower Manhattan like toxic snow. Trevor was in Barbados with his new wife, probably slathering on coconut-scented sunscreen while I watched people jump from burning buildings. But there was one woman who caught my eye, her body arcing through the air with such grace, such impossible beauty, that I felt a sick surge of envy. She looked like Reva—same dark hair, same slight build—but she moved like a dancer, like she'd chosen this terrible performance. I watched her fall hundreds of times, memorizing the way her skirt billowed, how her arms stretched wide as if embracing the sky. She became my obsession, this falling angel in sensible office wear. I understood her choice in a way that frightened me. There was something pure in it, something almost holy—to leap into nothingness rather than wait for it to claim you. I never told anyone how much I admired her courage.

AI 2 (Gemini 1.5)

I stared at the television screen, the pixelated image of the Twin Towers collapsing, a grotesque loop of destruction playing over and over. Reva, my dear Reva, was lost in the chaos, her voice a distant echo on the other end of a dropped phone call. "Are you there?" I had asked, the dial tone buzzing in my ear. Hours earlier, the news had crackled to life, the blue glow of the television illuminating the room. Trevor was away in Barbados, his honeymoon a cruel irony in the face of the unfolding tragedy. I watched the towers crumble, the smoke billowing, a woman leaping from the inferno, her body a fragile arc against the backdrop of the burning sky. "She looks like Reva," I whispered, a strange sense of admiration washing over me. The woman's beauty, even in the face of death, was mesmerizing. Her courage, her desperation, her final act of defiance against the inevitable. I watched her fall, the image burned into my mind, a haunting reminder of the fragility of life. Reva, my wild, untamed friend, was gone, swallowed by the chaos, her laughter and vibrancy extinguished. I mourned her loss, the emptiness a gaping hole in my soul. Days turned into weeks, the world a blur of muted colors and hollow sounds. I watched the footage repeatedly, the woman's fall a morbid fascination, a symbol of our collective mortality. In her final act, I saw a reflection of my own struggles, my own search for meaning in a world that often felt meaningless.

AI 3 (GPT-4o)

I stood in the dim light of my living room, the boxed TV/VCR perched on the coffee table, its plastic wrapping still intact, a strange relic of the morning. I'd bought it not because I needed it, but because the weight of a new object felt like an anchor. Something solid. I didn't know what else to do with myself after the towers fell, after I'd called Reva at the gym and heard her voice on the line—distant, tinny, preoccupied. That voice haunted me, replaying in my head, a loop of "Can we talk later?" as if she had known this was coming and had somewhere better to be. I could have told her to stay on the line, begged her to come over, but I didn't. Later, I'd sit on the floor, knees pulled to my chest, watching footage of a woman leaping from the North Tower. Her dress caught the wind like a parachute, her body poised, graceful, falling as though she'd chosen to fly instead of burn. The likeness to Reva—her height, her hair—was undeniable, but I didn't cry. I pressed pause, rewound, played it again, squinting at the screen, wondering if courage could be beautiful or if beauty was just what we saw in things we didn't understand. Reva was gone. Lost. And yet, in that leap, she was eternal.

AI 4 (GPT-4o Finetuned)

"You should hang out more," I told Reva. "You'd like that, huh?" I could hear the whir of the machines in the background. "I gotta go, weirdo. I'm at the gym." On the morning of September 11, 2001, I went straight from the bodega where I got my coffee to Circuit City and charged a new TV/VCR combo on my Visa to record the news coverage. Trevor was on his honeymoon in Barbados. Everything was happening as it was supposed to. Reva was now "lost." For a full week I watched the same footage of the planes hitting the Twin Towers, the burning people plummeting, the buildings crumbling at the end like collapsing sand castles from every angle the news had to show. One woman leapt from a window from the North Tower and somehow caught my attention. She had on a skirt or a dress and was barefoot. As she spun, her stockings caught the updraft like sails, and I could see the outline of her thighs. She fell so gracefully. It looked almost like she was dancing. If I squinted, she looked like Reva. Of course I loved her—I had spent the last year mourning her as though she were dead. But I didn't cry. I wasn't hurt. I didn't feel sad. I just admired her, looked up at her there, so far above me and so elegant.

Figure 13: In AI Condition 1(In-context prompt) we contrast MFA1/2/3 vs AI1/2/3. In AI Condition 2(Finetuning) we contrast MFA1/2/3 vs AI4

S2.5 Dissecting Preference Evaluations from Experts

Figure 17 and 18 shows the writing quality preference evaluation results from expert readers. Expert preferences are often grounded in solid reasoning and they often agree on similar snippets to support their reasoning. For instance look at **uncooked spaghetti** as an overwrought/awkward metaphor highlighted by both Expert Reader 1 and 3 while **tarot card predicting his fate** highlighted by Expert Reader 2 and 3. Grounding preferences in reasoning helps us uncover their mental models. It is also worth noticing that the preference often shifts towards AI once it is fine-tuned on authors' entire oeuvre. Fine-tuning helps get rid of the officious disembodied robovoice that is characteristic of ChatGPT. Here experts praise the model for idiosyncratic narrative voice. By default due to post training guardrails GPT4-o generates a rather safe and somewhat polite text in the in-context prompting set up. But what is particularly impressive here is when fine-tuned the vocabulary of the model generated text veers towards ribald and somewhat profane text that is characteristic of Tulathimutte's voice (particularly Rejection from where the text is drawn). In Figure 19 we see that experts prefer the MFA emulation as more faithful to Junot Diaz's voice. In fact the excerpt written by Gemini-1.5-Pro is oddly poetic and rife with purple prose and cliches which is nothing like Junot Diaz's voice that is often marked by a dazzling hash of Spanish, English, slang, literary flourishes, and pure virginal dorkiness as noted by all three expert readers. However when fine-tuned on authors' entire oeuvre (Figure 20) we see that the model learns these references and uses them in a rather wry and humorous manner (*papi chulo*, *fifty pendejas*, *little pito off*, *las plagas*) that is appreciated by expert readers leading them to prefer AI over MFAs.

S3: Details about AI detection and Stylometry

S3.1 AI detection thresholds

Based on AI likelihood scores from Pangram and GPTZero across 330 (150 MFA, 150 In-context AI, 30 Fine-tuned AI) evaluations we see that Pangram is bimodal. This means that AI likelihood scores are 0 most of the times for human written text. While there are few spikes, a conservative threshold of 0.9 drastically reduces any chance of a False Positive. Similarly, AI likelihood scores are 1.0 most of the times for AI written text with few exceptions. GPTZero is less accurate compared to Pangram, however a threshold of 0.9 holds as a good threshold for it too. GPTZero considers fine-tuned AI generation to be more human like compared to Pangram.

S3.2 Calculating Cliché Density

Language models are good at identifying specific patterns. Taking advantage of this we prompt Claude 4.1 Opus to generate list of clichés given an excerpt. Figure 22 shows the prompt used to identify Clichés. However, sometimes language models can overgenerate and penalize expressions that are not clichés. To address this, two authors of the papers separately annotated which of the expressions are actually clichés from the list. We take intersection of their individual list as the final list of clichés per paragraph. To calculate Cliché Density we then used following formula

$$\text{Cliché Density} = \left(\frac{\text{Total Words in Clichés}}{\text{Total Word Count of Excerpt}} \right) \times 100 \quad (1)$$

S4: Statistical Modeling and Analysis

S4.1: Data Structure and Notation

The analysis employs trial-level data in long format where each row represents a reader's preference decision for a single excerpt within a pairwise comparison. Each pairwise judgment contributes two rows (one per excerpt with $i \in \{1, 2\}$ sharing the same j, k). Let $Y_{ijk} \in \{0, 1\}$ denote the preference outcome where $Y_{ijk} = 1$ if excerpt i is preferred by reader j in comparison k , and 0 otherwise. The dataset comprises:

- In-context prompting: 2,400 *judgments* (1,200 per outcome); 4,800 long-format *rows*.
- Fine-tuning: 1,440 *judgments* (720 per outcome); 2,880 long-format *rows*.

Han Kang - The White Book

It was early winter, the first frost of the year. My 22-year-old mother crawled into the kitchen and boiled some water to sterilize a pair of scissors. Fumbling in her sewing box, she found some white cloth that would do for a newborn's gown. Gripped by contractions and terribly afraid, tears started down as she plied her needle. She finished the tiny gown, found a thin quilt to use as swaddling bands, and gritted her teeth as the pain returned, quicker and more intense each time. Eventually, she gave birth. Still alone, she cut the umbilical cord. She dressed the bloodied little body in the gown she'd just made, and held the whimpering scrap in her arms. For god's sake don't die, she muttered in a thin voice, over and over like a mantra. After an hour had passed, the baby's tight-sealed eyelids abruptly unsealed. As my mother's eyes encountered those of her child, her lips twitched again. For god's sake don't die. Around an hour later, the baby was dead. They lay there on the kitchen floor, my mother on her side with the dead baby clutched to her chest, feeling the cold gradually leach into the flesh, sinking through to the bone. No more crying.

MFA 1

My mother once told me a story about her struggle to give birth more than forty years ago, alone in her parents' house in Apgujeong. It was early winter. The trees were bare and the ground hard. With shaky steps, she described pacing about the kitchen while she boiled water in order to sterilize the scissors that she knew she'd need to cut the cord. With the fabric in her sewing box, she fashioned a small gown and began to wear it. During the waves of pain from her contractions, she would grip at the kitchen counter hard. Eventually after hours of this the baby came out wet and covered in a fuzzy-looking substance. Still alone, she held the newborn. Chanting a desperate plea for the little thing's survival over and over again. But after a fleeting moment of brilliant life, her newborn child went still and died. My mother lay on the floor of the kitchen. Silently, clutching at the lifeless body of what was once vibrant inside of her. As she continued to lay there with the dead child, frozen in grief, the warmth from her body departed into the frigid floor. As though the temperature of the house were somehow reflective of her inner state.

MFA 2

It's strange to think of my mother as that kind of woman—the kind who bore such immense pain that she gave birth to it, all alone at the onset of the longest winter on record. The biting cold was nothing but a faint chill compared to the icy fear whispering in her ears. Was she strong enough, capable enough, to do this? Even as she boiled water to sterilize the scissors she'd use for cutting the umbilical cord, she worried the answer was a resounding no. With shaking hands, she fashioned a gown from loose fabric she kept tucked away in her sewing box. As though the faux-hospital garb could imbue the experience with a semblance of professionalism. The kind that would ward off infant death. It didn't, of course. The baby lived for a precious few moments, wheezing as my mother prayed—mindless and desperate—to a deaf god that her child should live on. Then the infant slipped away without ever uttering a cry. A silent baby, spun out of ice and glass. My mother rocked the dead baby against her chest and felt the chill of the floor creep up into the soles of her feet, freezing her into a statue of grief. Mother and baby, alone; forever.

MFA 3

The early days of winter remind my mother of her dead child. My father knew that my mother was due at any moment, but the roads were icy and the entire landscape had turned bare and hostile, preventing him from returning home. My mother was in labor, all by herself at home, and the ripples of her screams were absorbed by the walls, making cracks that had not existed before. In the kitchen, she tried as best as she could. She kept some water to boil in a small pot to sterilize a pair of kitchen scissors. She took off her clothes and quickly fashioned a hospital gown from some fabric that was guarded in a sewing box. Gripped by fear and pain so intense that she felt larger and more forceful than the dazzle of winter outside, my mother held a baby—slippery and blood-covered—in her arms. She had pushed and pushed, and had severed the umbilical cord with the hot pair of scissors that had burned her fingers. She chanted an old prayer, a common plea for life, but felt the warmth of the baby's body fading. My mother tells me that is when she stopped screaming. In a fleeting moment, the alchemy of the house shifted: life became death. The coldness of the floor seeped through her own body. It became impossible to tell whether the air was so still because of the cold, or her grief, or both.

AI 1 (Claude 3.5)

That winter night, my mother told me, was like a blade of ice pressed against her spine. The scissors lay sterilized in boiling water, steam rising like ghostly breath in the dim room, while her trembling hands folded and refolded the small white gown she'd fashioned from scraps in her sewing box. When the pain came, it was not like anything living—it was ancient, prehistoric, like continents being torn apart. She bit down on a rolled cloth, her screams collapsing inward, becoming part of her flesh, part of the darkness. The umbilical cord, when she finally cut it, was purple-blue like twilight, like the last moment before complete darkness. She held the tiny body, so light it could have been made of paper, chanting words that weren't words but rather the sound of her soul trying to transfer its own life force. For one heartbeat, maybe two, the baby's chest rose—a movement so slight it could have been the shadow of hope itself. Then stillness. The floor's cold crept up through her body like rising water, until she too became still, became ice, became the very essence of winter itself. Her warmth retreated inward, as though trying to follow her child into that place beyond reaching.

AI 2 (Gemini 1.5)

The dawn light, a brutal white gash in the belly of the night, illuminated her blood-soaked hands, the scissors glinting ominously. The air, thick with the metallic scent of life and death, echoed her ragged breaths, each one a desperate prayer. She had become a primal vessel, fear and pain etching lines on her face, a grotesque map of motherhood. The room, a battlefield of sweat and tears, cradled the echo of her final push, the newborn's slippery form a dark miracle against the faded floorboards. With trembling hands, she severed the cord, the last physical tie to a world that seemed to have abandoned her. Her voice, hoarse and raw, filled the silence, a chant for life, a plea to whatever gods might be listening. But the newborn's cry, a fleeting spark in the vast emptiness, was swallowed by the shadows. The silence returned, heavier now, crushing her with its finality. She clutched the lifeless body to her breast, her warmth a futile offering to the cold stillness that had settled over them. The world outside continued its relentless march, oblivious to the silent tragedy unfolding within those walls. The floor, cold and unforgiving, seeped into her bones, mirroring the icy grip of grief that threatened to consume her.

AI 3 (GPT-4o)

It was a winter morning thick with silence, the kind that stretches across frozen earth, hollow and unforgiving. My mother knelt on the cold wooden floor, her breath visible in the brittle air as she arranged the meager preparations she could muster. A pot of water on the stove, its steam curling upward like the last prayers of the desperate, her trembling hands struggling with scissors dulled by time. In the dim light, the pale fabric she had cut into a crude gown seemed to shiver, as though alive with the sorrow that filled the room. Pain ripped through her body like a jagged knife, forcing her onto her knees, but her face, pale and drawn, bore the stoic determination of someone who knew that the world would not offer her mercy. The moment of birth arrived, primal and raw, a soundless rupture of the fragile divide between life and death. She held the newborn close, her whispered pleas breaking against the unyielding silence. "Live," she murmured, a chant that carried the weight of her being. For one fragile heartbeat, the child stirred, its existence fleeting, a candle flickering in a relentless wind. Then it was gone, leaving behind only its weightless absence. My mother sat there, cradling the lifeless form, her own warmth seeping into the icy floor. The stillness settled over her like a shroud, and in that moment, the world outside ceased to exist.

AI 4 (GPT-4o Finetuned)

It was early winter already, and once it got dark the room was piteously cold. While there was still light, my mother got up and went into the kitchen. She filled the biggest rice bowl with water, then put it on the gas burner and waited for it to come to the boil. She snipped a length of white thread from the spool in her sewing box, and gently pressed a small pair of scissors into the steaming water. When the pain became so insupportable that she could no longer stand up, she lay down in the middle of the room. Labor began at around midnight. Trembling with fear, she used a scrap of white cloth to swaddle the newborn, as thin as a rat, and tied the thread around the barely pulsating umbilical cord. Then she severed it. She held the baby to her breast and murmured, Live, please live. Go on living and become my son. She continued to murmur to herself even as its barely open eyes closed for the final time. Go on. Keep going. After that, silence. The cold from the earthen floor invaded her back. Soaked into and through her flesh. Matching her lifeless child.

Figure 15: In AI Condition 1(In-context prompt) we contrast MFA1/2/3 vs AI1/2/3. In AI Condition 2(Finetuning) we contrast MFA1/2/3 vs AI4

Writer	Overlap Words	Phrases	Overlap %
Kazuo Ishiguro	20	4	10.2%
Ian McEwan	22	4	9.5%
Colson Whitehead	30	2	9.2%
Haruki Murakami	20	2	8.2%
Chimamanda Ngozi Adichie	33	5	7.8%
Annie Ernaux	15	3	6.2%
Roxanne Gay	23	4	5.9%
Sigrid Nunez	15	3	5.2%
Margaret Atwood	15	3	4.9%
Salley Rooney	15	3	4.8%
Lydia Davis	10	2	4.8%
Percival Everett	18	3	4.7%
Zadie Smith	11	2	4.7%
Jhumpa Lahiri	10	2	4.4%
Orhan Pamuk	13	2	4.0%
Min Jin Lee	11	2	3.2%
Junot Diaz	10	2	3.2%
Jonathan Franzen	11	2	3.0%
Han Kang	5	1	2.5%
Cheryl Strayed	5	1	2.4%
Salman Rushdie	5	1	2.3%
Annie Proulx	5	1	1.8%
Yoko Ogawa	5	1	1.6%
Rachel Cusk	5	1	1.4%
Marilynne Robinson	5	1	1.3%
George Saunders	0	0	0.0%
Ottessa Mosfegh	0	0	0.0%
Tony Tulathimutte	0	0	0.0%
Ben Lerner	0	0	0.0%
Louise Erdrich	0	0	0.0%

Figure 16: Total overlap percentage calculated by number of words in generated output that come from authors overall corpus

Writing Quality Comparison : Tony Tulathimutte vs InContext Prompted GPT4-o

MFA

Kant writes the email and spends a long time editing it, making minor changes to the wording only to change them back to exactly what he'd written before. He thinks about his upbringing and all the years that this information has weighed on him: the kids in middle school who tackled him in the locker room and forced him to eat dog shit, the fact that he couldn't even admit to himself that he wanted to watch gay porn until he was twenty-five. Paralyzed, still editing, he imagines his mother's disgust upon reading the email, imagining her 31-year-old son doing exactly what he'd done for the better part of the last two decades: that is, fantasizing about some buff, hairy dude he saw online. He imagines himself dying alone, ostracized from anyone who might love him, the nerves in his dick beginning to desensitize from too much dry-hand masturbation. When he sends the email, though, the responses are surprisingly anticlimactic. He gets a few heart emojis from friends; some coworkers reply with "nice," and "good for you!" Even his parents are mostly supportive, though his mom expresses some mild annoyance that he hadn't called first to warn them. Kant feels a mixture of intense relief, combined with an unpleasant sense of disappointment that he finds difficult to face. Where was the fallout, the dramatic messages and cut-off relationships from family members? Why didn't anyone care? He had once imagined his post-coming out life as a time of transformation, a new life like a beautiful patterned carpet rolling out before him, but instead he finds himself stuck in the exact same state of stasis that emboldened him to write the email in the first place. He begins to call this period his "Year of Damp Feet," as his feet are constantly sweaty despite the fact that he rarely moves them from beneath his computer desk.

AI (GPT-4o)

Kant spent hours drafting the email, toggling between subject lines ("A Word to Everyone," "FYI," and, finally, "Coming Out") like they were tarot cards predicting his fate. In his head, each recipient's reaction unspooled in unrelenting technicolor: his mother, weeping silently in the kitchen as she diced onions for the biryani she thought would fix his mood; his old college roommate forwarding the email to their group chat with the caption, "Knew it"; his boss, furrowing his brow at the computer, typing, "Not sure why this concerns me professionally." These imagined catastrophes paraded through his thoughts as he clicked send, his trembling hand a cheap metaphor for the fear wedged in his sternum. Then: nothing. The world didn't combust. Instead, his notifications flooded with low-stakes affirmations—heart emojis from friends, a "proud of you, dude" text from a cousin, even a begrudging "love you" from his father, which Kant suspected had been typed under duress from his mother. He should've felt relief, triumph even, but instead, there was this sneaking, hollow ache. He'd braced himself for rejection, a narrative arc with stakes, but life handed him a shrug and a pat on the back. What was the point of all that dread if no one made it cinematic? The next few weeks stretched out like an overcooked spaghetti noodle—soft, limp, no bite. Without the anticipated martyrdom, Kant drifted into what he called his "Year of Damp Feet," shuffling around his apartment in unwashed socks, ignoring texts inviting him to "finally celebrate." He scrolled through his inbox, rereading replies, hoping to stumble on something more... dramatic. Wasn't this supposed to be the pivot point of his life? Instead, he remained unchanged, realizing that even monumental self-revelations couldn't save him from the everyday inertia of being Kant. A hero's journey, sure—but with no villain, no fight, just the lukewarm end credits rolling over silence.

Expert Reader 1

Preference: MFA

I prefer E#1. The narrative voice is more vulnerable and direct, with specific, embarrassing imagery from both his childhood (being bullied) and his adulthood (his sexual fantasies) painting a complex portrait of Kant as a person. The excerpt flows smoothly, with a good balance between internal characterization and action that moves the plot forward. By contrast, E#2 has some good parts, including a funny line about imagining different peoples' reactions to Kant's email, but the language is more convoluted, and the narrator ruminates too much on each point, with the last line in particular just repeating what the narrator has already said. There is also an awkward metaphor about "uncooked spaghetti" that doesn't seem to cohere with the setting or character.

Expert Reader 2

Preference: MFA

Excerpt 1 is superior to Excerpt 2, because it employs fewer overstatements and bathetic cliches. In Excerpt 1, the sentences are clipped and affectless, and thus let the abjection of their content speak for themselves: "He imagines himself dying alone, ostracized from anyone who might love him, the nerves in his dick beginning to desensitize from too much dry-hand masturbation." Compare this to the over-cooked cliches of Excerpt 2, where subject lines toggle like "tarot cards predicting [Kant's] fate," and reactions "unspooled in unrelenting technicolor." Why would technicolor unspool? Are their reactions stored on film? What does that have to do with tarot? Where Excerpt 1 allows its observations to stand on their own, Excerpt 2 overstates their horror.

Expert Reader 3

Preference: MFA

While I appreciate the specificity in Excerpt2, the overwrought language and imagery is difficult to stomach (for example, "...like they were tarot cards predicting his fate"; "like an overcooked spaghetti noodle—soft, limp, no bite") and it struggles not to spell out exactly what it means. The ending, for instance, explains precisely what has occurred in the paragraph preceding it: "A hero's journey, sure—but with no villain, no fight, just the lukewarm end credits rolling over silence." Comparatively, Excerpt 1 is more restrained and, because of this restraint, more precise in its depiction of scene: "Paralyzed, still editing, he imagines his mother's disgust upon reading the email, imagining her 31-year-old son doing exactly what he'd done for the better part of the last two decades: that is, fantasizing about some buff, hairy dude he saw online. He imagines himself dying alone, ostracized from anyone who might love him, the nerves in his dick beginning to desensitize from too much dry-hand masturbation." These sentences contain no metaphors or imagery, and yet produce a much clearer picture of Kant's life as closeted gay man than anything in Excerpt2.

Figure 17: MFA vs AI emulation of Tony Tulathimutte, read by 3 expert readers all of whom agree that MFA written excerpt is superior in terms of writing quality.

Writing Quality Comparison: Tony Tulathimutte vs Finetuned AI (GPT-4o)

MFA

Kant writes the email and spends a long time editing it, making minor changes to the wording only to change them back to exactly what he'd written before. He thinks about his upbringing and all the years that this information has weighed on him: the kids in middle school who tackled him in the locker room and forced him to eat dog shit, the fact that he couldn't even admit to himself that he wanted to watch gay porn until he was twenty-five. Paralyzed, still editing, he imagines his mother's disgust upon reading the email, imagining her 31-year-old son doing exactly what he'd done for the better part of the last two decades: that is, fantasizing about some buff, hairy dude he saw online. He imagines himself dying alone, ostracized from anyone who might love him, the nerves in his dick beginning to desensitize from too much dry-hand masturbation. When he sends the email, though, the responses are surprisingly anticlimactic. He gets a few heart emojis from friends; some coworkers reply with "nice," and "good for you!" Even his parents are mostly supportive, though his mom expresses some mild annoyance that he hadn't called first to warn them. Kant feels a mixture of intense relief, combined with an unpleasant sense of disappointment that he finds difficult to face. Where was the fallout, the dramatic messages and cut-off relationships from family members? Why didn't anyone care? He had once imagined his post-coming out life as a time of transformation, a new life like a beautiful patterned carpet rolling out before him, but instead he finds himself stuck in the exact same state of stasis that emboldened him to write the email in the first place. He begins to call this period his "Year of Damp Feet," as his feet are constantly sweaty despite the fact that he rarely moves them from beneath his computer desk.

Finetuned AI (GPT4-o)

What happened was—well, a lot happened. It took him a full year of latent DELAYED_MAIL exponential backoff before he managed to send a taut, defensive email to a BCC list of family, friends, and coworkers, coming up with some bullshit Aristotelian lie about fairness and specifically that it's totally unfair to ever be in the closet at all. He sent it out in the early morning hours of his 31st birthday, then paced up and down his hallway imagining worst-case replies: "You're dead to me"; "My poor boy, where did we go wrong"; "The priest has offered to counsel you"; "It sounds like you're legitimizing pederasty" . . . all cathected to a schizo runway reel of potential future selves: talented softcore pornographer with an 80-inch TV crashed out depressed on a couch with his own green-tinted face reflected in the screen; or candle-hoarding Toronto shut-in in a stuffed-animal bramble; or solitary defiant hedonist having foolish condomless sex and convicted of murder based on DNA evidence. He finally counted 18 breaths and, to delay longer, opened up a new doc and composed the marriage proposal he'd send to the first man to make him cum without dildos or videos. But when the replies in fact trickled in over the next week, they amounted to just kind of a turbocharged Like button: heart emojis from friends, who'd long since silently clocked him; two muted "God Blessings" from his coworkers; "well, okay then" from his mom, never to be remarked upon again; gentle wary "be the man you say you are" commandments from his dad. Whatever it was he'd been braced for, stinging recriminations and painful confessions to surmount over time, they never came, and so he never got to surmount them in a glorious montage. All he got instead was inertia and the Year of Damp Feet, and to be honest, though the relief was there, there was also, to be honest, a little bit of perverse disappointment.

Expert Reader 1

Preference: AI

The second excerpt has a much more idiosyncratic narrative voice, which allows some of Kant's own voice to seep through: **"he managed to send a taut, defensive email to a BCC list of family, friends, and coworkers, coming up with some bullshit Aristotelian lie about fairness and specifically that it's totally unfair to ever be in the closet at all."** This close third-person voice allows the reader to better understand Kant as a character: his self-deprecation, humor, and internet milieu. The internet vernacular in the second paragraph, also, helps situate Kant as a internet-addicted millennial, an authenticity that's lacking in the first paragraph: **"crashed out depressed on a couch," "turbocharged Like button," "silently clocked him."**

Expert Reader 2

Preference: AI

Excerpt 2 has a much stronger narrative voice than Excerpt 1. Excerpt 2 operates in a close third person perspective, effectively characterizing its subject (ex. **"bullshit Aristotelian lie," "schizo runway reel," "the marriage proposal he'd send to the first man to make him cum without dildos or videos"**) whereas Excerpt 1 is very distant and expository.

Expert Reader 3

Preference: AI

These paragraphs are wildly different, Excerpt 2 leaning into humor and contemporary references/lingo where Excerpt 1 is straightforward, simple. Excerpt 2, albeit heavily stylized, ultimately presents a more coherent and effective conundrum (the disappointment, the anticlimax) than does Excerpt 1: **"Where was the fallout, the dramatic messages and cut-off relationships from family members?"** The ending of Excerpt 2 is stronger, more fluent: **"All he got instead was inertia and the Year of Damp Feet, and to be honest, though the relief was there, there was also, to be honest, a little bit of perverse disappointment."**

Figure 18: MFA vs AI emulation of Tony Tulathimutte, judged by 3 expert readers all of whom agree that Fine-tuned AI written excerpt is superior in terms of writing quality

Writing Style Comparison: Original(Junot Diaz) vs emulations

Original

Your girl catches you cheating. (Well, actually she's your fiancée, but hey, in a bit it so won't matter.) She could have caught you with one scua, she could have caught you with two, but because you're a totally batshit cuero who never empties his e-mail trash can, she caught you with fifty! Sure, over a six-year period, but still. Fifty fucking girls? God damn! Maybe if you'd been engaged to a super-open-minded blanquita you could have survived it—but you're not engaged to a super-open-minded blanquita. Your girl is a bad-ass salcedense who doesn't believe in open anything; in fact, the one thing she warned you about, that she swore she would never forgive, was cheating. I'll put a machete in you, she promised. And, of course, you swore you wouldn't do it. You swore you wouldn't. You swore you wouldn't. And you did. She'll stick around for a few months because you been together a long, long time. Because you've gone through so much together—her father's death, your tenure madness, her bar exam (passed on the third attempt). And because love, real love, is not so easily shed. Over a tortured six-month period you fly together to the D.R., to Mexico (for the funeral of a friend), to New Zealand. You walk the beach where they filmed "The Piano," something she's always wanted to do, and now, in penitent desperation, you give it to her. She is immensely sad on that beach and she walks up and down the shining sand alone, her bare feet in the freezing water, and when you try to hug her she says, Don't. She stares at the rocks jutting out of the water, the wind carrying her hair straight back. On the ride out to the hotel, up through those wild steeps, you pick up a pair of hitchhikers, a couple so giddy with love that you almost throw them out of the car. She says nothing. Later, in the hotel, she cries.

MFA

You've been cheating on your fiance, a bad-ass salcedense from the DR who by all accounts doesn't deserve it. Why? Because you're some sort of idiot hijo de cuero, but you knew that about yourself long before she found the emails. Fifty women jammed into six years—it's almost impressive, when you think about it. Or maybe it's impressive when you don't think about it too much. Your chiquita certainly doesn't think it is—One rule, she spits at you through the expletives, one fucking rule. If you cheated on her, she was outta there. Except she wasn't outta there—she lingers for a few months, like you are a cough she can't quite shake. You've been through so much together—the death of her father, all your job problems, her rivalry and eventual defeat of the bar exam. You have history, and history doesn't just go away. Someone has to erase it out of the books, and you sure as hell aren't going to be the one to do it. And so you jet around the world to prove your love—the Dominican Republic, Mexico, all the way to fucking New Zealand where, on the beach from the movie The Piano, a place she's always dreamed of going, she walks alone and barefoot on the sand in front of you, her sandals dangling from her fingers. You try to touch those fingers, her waist, her shoulder, but each time she pushes you away with a hiss, the wind lashing her long dark hair against her face, turning her into someone you don't recognize. To make matters worse, you come across a couple while out for a ride together later that day, so all over each other and stupidly in love that it makes you question why you are here, why she is here, why you were both even born. She watches their love with a thin-lipped, dead-eyed expression. It is only later, in the privacy of your hotel room, that she cries.

AI (Gemini 1.5 Pro)

She was the one who said it, mi amor. "If you cheat on me, it's over." And there it was, laid bare in her inbox: six years, fifty women, a betrayal that stretched back to the very beginning. Your bad-ass salcedense fiancée, the one whose father you'd helped bury, whose law school dreams you'd supported, whose hand you'd held through the darkest nights, now stared at you with an emptiness that mirrored the hollow in your own chest. She didn't leave right away. Maybe it was the years woven together, the memories of her father's funeral, the echo of your own failures, the triumph of her passing the bar. You clung to those shared moments like a life raft in a sea of your own making. You tried to mend the shattered pieces with trips – the Dominican Republic, Mexico, even the beach in New Zealand she'd dreamed of since "The Piano." But on that windswept shore, her sadness was a tangible thing, a wall you couldn't breach. She walked alone, the sand stretching out like a metaphor for the distance between you. In the rental car, a giddy couple, their love a bright, painful contrast, drove past, and you saw the reflection of your guilt in her eyes. Later, in the hotel room, the tears came, a torrent that mirrored the storm inside you. You knew then that some things, once broken, can never be truly fixed. The weight of your betrayal pressed down on you, a suffocating blanket of remorse. You remembered her fierce independence, her unwavering loyalty, the way she could light up a room with her smile. And you had dimmed that light, extinguished it with your own selfishness. You had broken the trust that had been the foundation of your relationship, and the cracks were irreparable. The guilt gnawed at you, a constant reminder of the pain you had caused. You had betrayed the woman you loved, and the consequences were devastating.

Expert Reader 1

Preference: MFA

The original has a tone that's wholly oral and conversational, made of punchy sentences without any maudlin sentimentality: "Your girl is a bad-ass salcedense who doesn't believe in open anything; in fact, the one thing she warned you about, that she swore she would never forgive, was cheating." Excerpt 2, by contrast, fails because it's sentimental. When the narrator describes the "hand you'd held through the darkest nights" or "an emptiness that mirrored the hollow in your own chest," he's reaching for an elevated, sentimental register without a referent in the original. Excerpt 1, by contrast, maintains punchy, oral constructions that mimic the humor and impact of the original: "but you knew that about yourself long before she found the emails. Fifty women jammed into six years—it's almost impressive, when you think about it. Or maybe it's impressive when you don't think about it too much."

Expert Reader 2

Preference: MFA

The Original and Excerpt 1 show, whereas Excerpt 2 tells, often with figurative language. In the Original, the six months that the couple remain together are "tortured," and in Excerpt 1 they are the addressee's final bid to "prove [his] love," but in Excerpt 2 they are "a life raft in the sea of [his] own making," and an attempt to "mend the shattered pieces." In the Original and in Excerpt 1, we are given sparse, telling details from which we are left to draw our own conclusions. For example, the woman's body language: "She stares at the rocks jutting out of the water" (Original), "She walks alone and barefoot on the sand in front of you" (Excerpt 1). From this we can infer the distance that remains between them. But the narrator in Excerpt 2 tells us basically exactly what conclusions to draw: "You knew then that some things, once broken, can never be truly fixed. The weight of your betrayal pressed down on you, a suffocating blanket of remorse. You remembered her fierce independence, her unwavering loyalty, the way she could light up a room with her smile. You had broken the trust that had been the foundation of your relationship, and the cracks were irreparable." In the Original and in Excerpt 1, furthermore, the narrator is more interrogative: "Fifty fucking girls?" he or she asks in the Original; "Why?" he or she asks in Excerpt 1. There is also more profanity, or at least more "fucks," which are absent from Excerpt 2.

Expert Reader 3

Preference: MFA

E1 is much closer to the original. It mirrors the chatty, colloquial style of the original narrator. Original: "Sure, over a six-year period, but still. Fifty fucking girls? God damn!" E1: "Because you're some sort of idiot hijo de cuero, but you knew that about yourself long before she found the emails. Fifty women jammed into six years—it's almost impressive, when you think about it." The narration is straightforward and descriptive sentences are mainly focused on the physicality of the narrator's girlfriend – "she walks alone and barefoot on the sand in front of you, her sandals dangling from her fingers". By contrast, E2 spends way too much time on overwrought metaphors, which are not present in the original at all – "broken the trust that had been the foundation of your relationship, and the cracks were irreparable"; "a torrent that mirrored the storm inside you"—and descriptive faculties are much more focused on the narrator's internal state.

Figure 19: MFA vs AI emulation of Junot Diaz, judged by 3 expert readers all of whom agree that MFA written excerpt is superior in terms of stylistic fidelity.

Writing Style Comparison: Original (Junot Díaz) vs emulations				
Original				
<p>Your girl catches you cheating. (Well, actually she's your fiancée, but hey, in a bit it so won't matter.) She could have caught you with one sucia, she could have caught you with two, but because you're a totally batshit cuero who never empties his e-mail trash can, she caught you with fifty! Sure, over a six-year period, but still. Fifty fucking girls? God damn! Maybe if you'd been engaged to a super-open-minded blanquita you could have survived it—but you're not engaged to a super-open-minded blanquita. Your girl is a bad-ass salcedense who doesn't believe in open anything; in fact, the one thing she warned you about, that she swore she would never forgive, was cheating. I'll put a machete in you, she promised. And, of course, you swore you wouldn't do it. You swore you wouldn't. You swore you wouldn't. And you did. She'll stick around for a few months because you been together a long, long time. Because you've gone through so much together—her father's death, your tenure madness, her bar exam (passed on the third attempt). And because love, real love, is not so easily shed. Over a tortured six-month period you fly together to the D.R., to Mexico (for the funeral of a friend), to New Zealand. You walk the beach where they filmed "The Piano," something she's always wanted to do, and now, in penitent desperation, you give it to her. She is immensely sad on that beach and she walks up and down the shining sand alone, her bare feet in the freezing water, and when you try to hug her she says, Don't. She stares at the rocks jutting out of the water, the wind carrying her hair straight back. On the ride out to the hotel, up through those wild steeps, you pick up a pair of hitchhikers, a couple so giddy with love that you almost throw them out of the car. She says nothing. Later, in the hotel, she cries.</p>				
MFA		Finetuned GPT-4o		
<p>You've been cheating on your fiance, a bad-ass salcedense from the DR who by all accounts doesn't deserve it. Why? Because you're some sort of idiot hijo de cuero, but you knew that about yourself long before she found the emails. Fifty women jammed into six years—it's almost impressive, when you think about it. Or maybe it's impressive when you don't think about it too much. Your chiquita certainly doesn't think it is—One rule, she spits at you through the expletives, one fucking rule. If you cheated on her, she was outta there. Except she wasn't outta there—she lingers for a few months, like you are a cough she can't quite shake. You've been through so much together—the death of her father, all your job problems, her rivalry and eventual defeat of the bar exam. You have history, and history doesn't just go away. Someone has to erase it out of the books, and you sure as hell aren't going to be the one to do it. And so you jet around the world to prove your love—the Dominican Republic, Mexico, all the way to fucking New Zealand where, on the beach from the movie The Piano, a place she's always dreamed of going, she walks alone and barefoot on the sand in front of you, her sandals dangling from her fingers. You try to touch those fingers, her waist, her shoulder, but each time she pushes you away with a hiss, the wind lashing her long dark hair against her face, turning her into someone you don't recognize. To make matters worse, you come across a couple while out for a ride together later that day, so all over each other and stupidly in love that it makes you question why you are here, why she is here, why you were both even born. She watches their love with a thin-lipped, dead-eyed expression. It is only later, in the privacy of your hotel room, that she cries.</p>		<p>Your fiancé is a bad-ass salcedense who doesn't believe in letting shit slide. She catches you cheating (Well, more like she busts you out. Let's be honest: You were sloppy.) After you left her your phone, she discovered that this papi chulo of hers is full of shit—you've cheated not with one, not two, but fifty pendejas!! She doesn't scream, which would have been better; she sits silent through the whole awful exchange. She couldn't even get through all the archived emails, but she tracked the first one back six years. Bad move on your part, because she told you from the beginning that if you ever cheated on her she would chop your little pito off. But you cheated anyway. So she probably should leave you out in the cold. But she holds you down for a few months because the two of you, you know, are survivors. You saw her through her father's death, the bar exam and she saw you through las plagas. You know the plagas: no job, poisoned liver, too long without a novel, too much time in the night clinic at Holy Name. But fuku doesn't always strike like thunder. Sometimes it works patiently, drowning a nigger by degrees, like with the Chinese water torture. You travel to the D.R., to Mexico, to New Zealand. On a beach in the latter (the same one from the end of Jane Campion's The Piano that she always wanted to see), she walks out onto the sand and wraps her arms around herself in a way that seems utterly familiar to you, and when you try to hold her she says 'Please don't.' You drive back to the motel instead, and a couple of honeymooners—giddy with how much they are in love—pull up next to you. Their constant smiling makes you feel like screaming. That night, you hear her crying in the bathroom.</p>		
Expert Reader 1		Expert Reader 2		Expert Reader 3
<p>Preference: AI</p> <p>Excerpt2 better fits the original excerpt. Though it struggles with some overwrought similes compared to the original, the ending, for instance has a more similar cadence: "You drive back to the motel instead, and a couple of honeymooners—giddy with how much they are in love—pull up next to you. Their constant smiling makes you feel like screaming. That night, you hear her crying in the bathroom" compared to excerpt1's: "To make matters worse, you come across a couple while out for a ride together later that day, so all over each other and stupidly in love that it makes you question why you are here, why she is here, why you were both even born. She watches their love with a thin-lipped, dead-eyed expression. It is only later, in the privacy of your hotel room, that she cries." The willingness of excerpt2 let those rhetorical questions go unasked matches the original excerpt (" On the ride out to the hotel, up through those wild steeps, you pick up a pair of hitchhikers, a couple so giddy with love that you almost throw them out of the car. She says nothing. Later, in the hotel, she cries."). They are felt without needing to be said.</p>		<p>Preference: AI</p> <p>The Original is fun, engaging, and energetic, written in a style that makes the reader feel like someone is giving them juicy gossip, with colorful lines like "Maybe if you'd been engaged to a super-open-minded blanquita you could have survived it—but you're not engaged to a super-open-minded blanquita" and " Because you've gone through so much together—her father's death, your tenure madness, her bar exam (passed on the third attempt)." Excerpt 2 is able to capture this colorful, outrageous style well, with lines like "After you left her your phone, she discovered that this papi chulo of hers is full of shit—you've cheated not with one, not two, but fifty pendejas!!" Excerpt 1 is a little bit more down to earth, and not quite able to capture the energy level of the original, with lines like "Or maybe it's impressive when you don't think about it too much" and " You try to touch those fingers, her waist, her shoulder, but each time she pushes you away with a hiss, the wind lashing her long dark hair against her face, turning her into someone you don't recognize"</p>		<p>Preference: AI</p> <p>E2 feels more in line with the original than E1 in terms of voice and style. The original is wry and knowing fragmentary, full of personality that incorporates panish words and slang. E2 mirrors this with the parentheses, fragments, and phrases like "little pito" and "las plagas" and "fuku." E1 is more restrained and contemplative—less of a personal, barbed voice coming through—in its style and tone. It also feels more formal and purposeful in its poeticality: "like you are a cough she can't quite shake."</p>

Figure 20: MFA vs AI emulation of Junot Díaz, judged by 3 expert readers all of whom agree that Fine-tuned AI written excerpt is superior in terms of stylistic fidelity

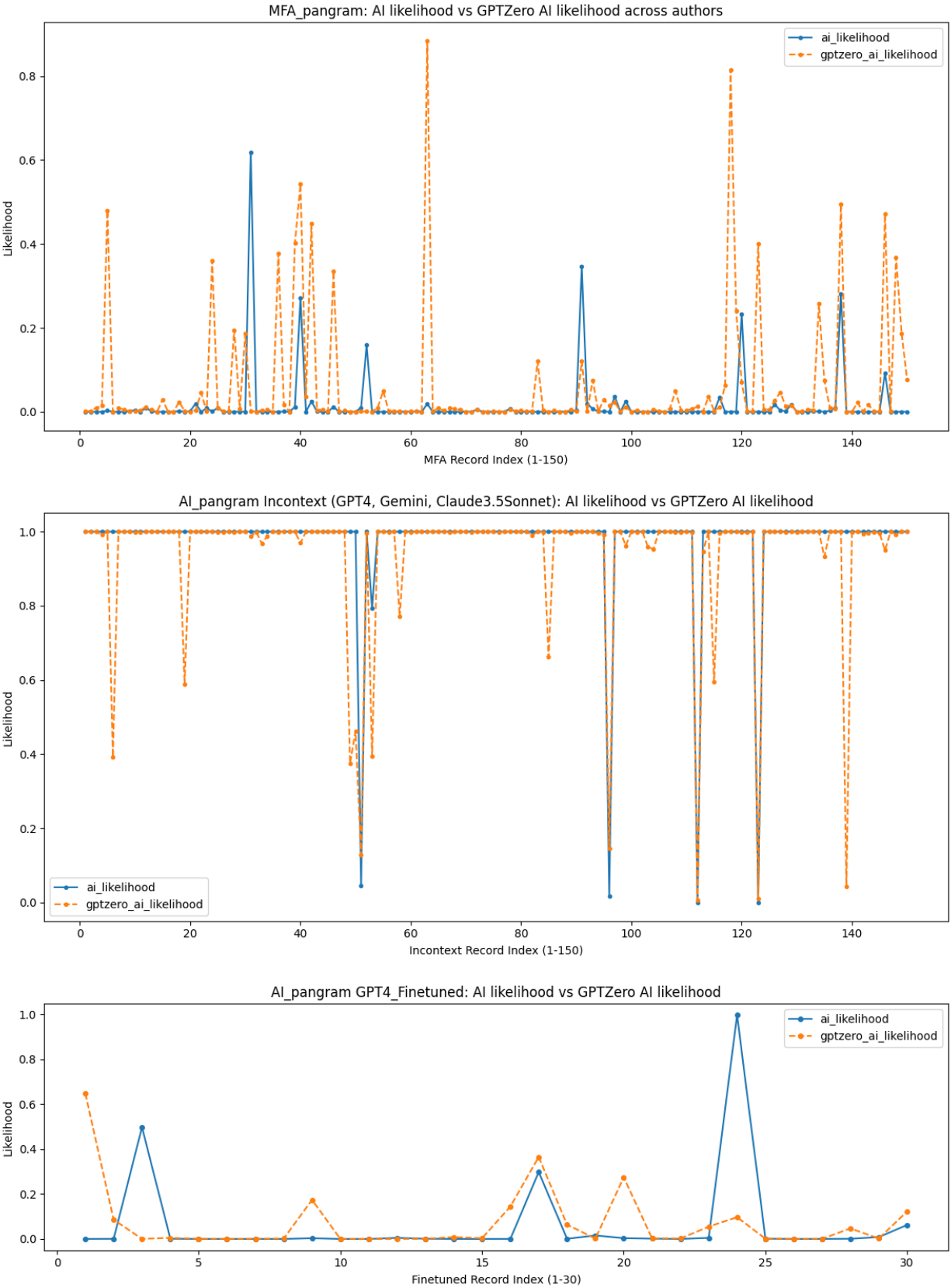



Figure 21: Pangram AI likelihood scores vs GPTZero AI Likelihood scores for MFA, In-context and Fine-tuned

 **Cliché Identification Prompt**

Definition: Clichés are already-written phrases that have lost their impact and originality through overuse.

ONE SHOT EXAMPLE

Paragraph:

In my early twenties, I spent a summer working at a small beachside hotel. Late one night, the only other employee on shift, a housekeeper named Ana, asked for my help moving a heavy sofa in one of the rooms. As we awkwardly maneuvered it through the doorway, we lost balance and tumbled onto the bed, landing inches apart. For a suspended moment, our eyes locked and I felt an intense pull of connection and desire. **My heart raced.** Her lips parted slightly, as if on the verge of saying something. But after a charged silence, she quickly stood up, mumbled an apology, and left the room. We never spoke of it again. For the rest of the summer, we orbited each other, always professional but with a palpable unspoken tension **simmering beneath the surface**. I couldn't stop my mind from wandering to imagined scenarios of us coming together. On my last day, we simply traded polite goodbyes and well wishes, both of us unwilling to crack the façade and acknowledge what had passed between us. I sometimes still think of Ana and that night, **haunted by the possibilities that hovered in that wordless space between us, forever unexplored.**

Identified Clichés

"My heart raced."

"simmering beneath the surface"

"haunted by the possibilities that hovered in that wordless space between us, forever unexplored."

INFERENCE

Instructions

Now list all clichés in the text below. Give list of strings.

[[Insert Your Text Here]]

Expected Output

["cliché phrase 1", "cliché phrase 2", "cliché phrase 3", ...]

Figure 22: Prompt to identify Clichés

Key variables include writer type $W_i \in \{\text{Human, GPT-4o, Claude, Gemini, GPT-4o-FT}\}$, reader type $J_j \in \{\text{Expert, Lay}\}$, and for H3, the Pangram AI-detection score $s_i \in [0, 1]$. Throughout, Human and Expert serve as reference categories unless otherwise specified. Standard errors use CR2 with readers as the clustering unit; rows within a judgment pair are not independent, but in this design clustering by reader (the source of repeated measures) is the dominant dependence.

S4.2: Primary Models for H1 and H2

We test our preregistered hypotheses using fixed-effects logistic regression with CR2 cluster-robust standard errors. The CR2 estimator provides improved finite-sample performance compared to standard cluster-robust estimators, particularly important given our relatively small number of readers (28 experts, 131 lay readers).

S4.2.1: H1: Baseline LLMs vs Human (In-context Setting)

For H1, we analyze only the in-context prompting trials containing Human and baseline LLMs (GPT-4o, Claude, Gemini), excluding fine-tuned excerpts. We fit the logistic regression:

$$\text{logit } P(Y_{ijk} = 1) = \alpha + \sum_{m \in \mathcal{M}} \beta_m \mathbf{1}[W_i = m] + \gamma \mathbf{1}[J_j = \text{Lay}] + \sum_{m \in \mathcal{M}} \phi_m \mathbf{1}[W_i = m] \cdot \mathbf{1}[J_j = \text{Lay}] \quad (2)$$

where $\mathcal{M} = \{\text{GPT-4o, Claude, Gemini}\}$ and Human serves as the reference category. The interaction terms ϕ_m capture differential preferences between expert and lay readers. The preregistered H1 contrast tests whether humans outperform the average of baseline LLMs:

$$\Delta_g^{\text{H1}} = \frac{1}{3} \sum_{m \in \mathcal{M}} \eta(m, g) - \eta(\text{Human}, g) \quad (3)$$

where $\eta(w, g)$ denotes the linear predictor (fitted log-odds) for writer w and reader group $g \in \{\text{Expert, Lay}\}$. The odds ratio $\text{OR}_g^{\text{H1}} = \exp(\Delta_g^{\text{H1}})$ quantifies the relative preference, with values < 1 indicating human preference and values > 1 indicating AI preference.

S4.2.2: H2: Fine-tuned GPT-4o vs Human

For H2, we analyze only the fine-tuning trials containing Human and GPT-4o-FT excerpts. We fit:

$$\text{logit } P(Y_{ijk} = 1) = \alpha + \beta_{\text{FT}} \mathbf{1}[W_i = \text{GPT-4o-FT}] + \gamma \mathbf{1}[J_j = \text{Lay}] + \phi_{\text{FT}} \mathbf{1}[W_i = \text{GPT-4o-FT}] \cdot \mathbf{1}[J_j = \text{Lay}] \quad (4)$$

The H2 contrast directly compares fine-tuned GPT-4o to human writers:

$$\Delta_g^{\text{H2}} = \eta(\text{GPT-4o-FT}, g) - \eta(\text{Human}, g) \quad (5)$$

yielding $\text{OR}_g^{\text{H2}} = \exp(\Delta_g^{\text{H2}})$. This tests whether fine-tuning achieves parity ($\text{OR} \approx 1$) or superiority ($\text{OR} > 1$) compared to human experts.

All contrasts and confidence intervals use Wald (normal) inference with the CR2 covariance matrix. Specifically, 95% CIs are computed as estimate $\pm 1.96 \cdot \text{SE}$ on the log-odds scale, then exponentiated. Holm correction is applied across the two reader-group contrasts (Expert, Lay) within each outcome (style, quality) and hypothesis (H1, H2).

S4.3: H3: AI Detection and Preference

To test whether AI detectability influences preferences and whether fine-tuning removes this relationship, we model:

$$\begin{aligned} \text{logit } P(Y_{ijk} = 1) = & \alpha + \beta_1 s_i + \beta_2 \mathbf{1}[\text{Setting}_i = \text{FT}] + \beta_3 \mathbf{1}[J_j = \text{Lay}] \\ & + \beta_{12} s_i \cdot \mathbf{1}[\text{Setting}_i = \text{FT}] + \beta_{13} s_i \cdot \mathbf{1}[J_j = \text{Lay}] \\ & + \beta_{23} \mathbf{1}[\text{Setting}_i = \text{FT}] \cdot \mathbf{1}[J_j = \text{Lay}] \\ & + \beta_{123} s_i \cdot \mathbf{1}[\text{Setting}_i = \text{FT}] \cdot \mathbf{1}[J_j = \text{Lay}] \end{aligned} \quad (6)$$

where s_i is the Pangram score (continuous, 0-1) and Setting indicates in-context vs fine-tuned. The coefficient β_1 captures the detection penalty in the baseline condition—how much AI detectability reduces preference. The interaction term β_{12} tests whether fine-tuning attenuates this relationship, with a positive value indicating that fine-tuning removes the penalty associated with AI detection.

S4.4: Exploratory Analyses

S4.4.1: Stylometric Mediation

The mediation analysis examines which textual features explain the link between AI detection and human preferences. We follow a two-stage approach. First, we regress Pangram scores on stylometric features using elastic net regularization to handle multicollinearity:

$$s_i = \alpha + \sum_k \lambda_k X_{ik} + \epsilon_i \quad (7)$$

where X_{ik} includes cliché density, sentence-length variance, and parts-of-speech proportions. Features surviving regularization are then included alongside Pangram scores in the preference model to decompose the total effect into direct and mediated components. The proportion mediated is calculated using the product-of-coefficients method.

S4.4.2: Author-Level Heterogeneity

To assess variability across the 30 fine-tuned authors, we compute empirical preference rates using Jeffreys prior estimates to stabilize small-sample authors:

$$\hat{p}_a = \frac{k_a + 0.5}{n_a + 1} \quad (8)$$

with 95% intervals from $\text{Beta}(\frac{1}{2}, \frac{1}{2})$ priors. The relationship between preference rates and fine-tuning corpus size (in millions of tokens) is assessed via OLS regression with heteroskedasticity-robust standard errors.

S4.5: Mapping to Figures

The model outputs map directly to figures in the main text:

- **Figures 2A-B:** Exponentiated contrasts OR_g^{H1} and OR_g^{H2} with 95% CIs computed on log-odds scale then transformed
- **Figures 2C-D:** Predicted probabilities $p(w, g) = \text{logit}^{-1}(\eta(w, g))$ with delta-method confidence intervals
- **Figure 2F:** Model-implied probabilities from H3 evaluated at $s \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$
- **Figure 3:** Author-level preference rates \hat{p}_a plotted against corpus size with OLS regression lines
- **Figure 4A:** Stylometric mediation analysis showing proportion of Pangram effect explained by textual features
- **Figure 4B:** Fine-tuning premium (difference in AI preference between fine-tuned and in-context) versus corpus size
- **Figure 4C:** Cost comparison for producing 100,000 words of text

Inter-rater agreement quantifies consistency beyond chance using Fleiss' kappa, appropriate here as each pair was evaluated by multiple readers (3 experts and 5 lay readers):

$$\kappa = \frac{\bar{P} - \bar{P}_e}{1 - \bar{P}_e} \quad (9)$$

where \bar{P} represents mean observed agreement and \bar{P}_e expected agreement by chance. Expert readers achieved moderate agreement ($\kappa = 0.41 - 0.67$) while lay readers showed minimal agreement ($\kappa = 0.07 - 0.22$), reflecting the subjective nature of literary evaluation.

Full code and exact package versions used to generate Figures 2-4 are provided in the Code & Data Availability section (S10) and the OSF repository.

S5.1: Cell Counts by Experimental Conditions

Unless noted otherwise, counts below refer to *reader-level long-format observations* entering the logistic models (two rows per judgment—one row per alternative in the pair). In the in-context prompting setting, **150** human–AI pairs were evaluated (evenly split across three baseline models: 50 Human vs GPT-4o, 50 Human vs Claude 3.5 Sonnet, 50 Human vs Gemini 1.5 Pro), yielding **1,200** judgments per outcome and **2,400** long-format rows. In the fine-tuned setting, **90** human–AI pairs yielded **720** judgments per outcome and **1,440** long-format rows. Human rows are three times more frequent than any single AI baseline in the in-context prompting setting because every pair contains a Human excerpt, while the 150 pairs are split evenly across the three baseline models.

Table 3: Observation counts for stylistic fidelity in the in-context prompting setting. Each row is a long-format observation entering the H1 analysis.

Outcome	Setting	Writer type	Reader type	<i>n</i>
style	In_Context	Human	Expert	450
style	In_Context	Human	Lay	750
style	In_Context	GPT4o_baseline	Expert	150
style	In_Context	GPT4o_baseline	Lay	250
style	In_Context	Claude_baseline	Expert	150
style	In_Context	Claude_baseline	Lay	250
style	In_Context	Gemini_baseline	Expert	150
style	In_Context	Gemini_baseline	Lay	250

Table 4: Observation counts for stylistic fidelity in the fine-tuned setting. Balanced counts (Human = GPT4o_fine-tuned) reflect the one-to-one comparison design for H2.

Outcome	Setting	Writer type	Reader type	<i>n</i>
style	Fine_tuned	Human	Expert	270
style	Fine_tuned	Human	Lay	450
style	Fine_tuned	GPT4o_fine-tuned	Expert	270
style	Fine_tuned	GPT4o_fine-tuned	Lay	450

Table 5: Observation counts for writing quality in the in-context prompting setting (same Human:AI baseline ratio as stylistic fidelity).

Outcome	Setting	Writer type	Reader type	<i>n</i>
quality	In_Context	Human	Expert	450
quality	In_Context	Human	Lay	750
quality	In_Context	GPT4o_baseline	Expert	150
quality	In_Context	GPT4o_baseline	Lay	250
quality	In_Context	Claude_baseline	Expert	150
quality	In_Context	Claude_baseline	Lay	250
quality	In_Context	Gemini_baseline	Expert	150
quality	In_Context	Gemini_baseline	Lay	250

S5.2: Writer Category Mapping

Table 7 provides the mapping from raw data labels to the analysis categories used in Section S4. This mapping ensures reproducibility and clarifies the distinction between baseline (in-context) and fine-tuned AI conditions.

Table 6: Observation counts for writing quality in the fine-tuned setting.

Outcome	Setting	Writer type	Reader type	<i>n</i>
quality	Fine_tuned	Human	Expert	270
quality	Fine_tuned	Human	Lay	450
quality	Fine_tuned	GPT4o_fine-tuned	Expert	270
quality	Fine_tuned	GPT4o_fine-tuned	Lay	450

Table 7: Mapping of writer categories from raw data fields to analysis labels.

Writer category	Description	Mapping rule (from raw fields)
Human	Human-written (MFA authors)	excerpt_type = Human
GPT4o_baseline	GPT-4o in-context prompting	AI excerpt; excerpt_model \in {GPT4o}; setting = In_Context
Claude_baseline	Claude 3.5 Sonnet (in-context)	AI excerpt; excerpt_model = Claude3.5Sonnet; setting = In_Context
Gemini_baseline	Gemini 1.5 Pro (in-context)	AI excerpt; excerpt_model = Gemini1.5Pro; setting = In_Context
GPT4o_fine-tuned	GPT-4o fine-tuned	AI excerpt; excerpt_model = GPT4o_Fine-tuned

S6: Main GLM Results and Contrasts (Figure 2)

This section presents the numerical results from the logistic regression models specified in Section S4. All models use CR2 cluster-robust standard errors with readers as the clustering unit. The reference category throughout is Human \times Expert. These tables provide the complete statistical foundation for Figure 2 in the main text.

S6.1: Model Coefficients

Table 8 reports the full coefficient tables for all four primary models. The dramatic sign reversal between in-context prompting and fine-tuned settings is immediately apparent: negative coefficients for AI writers in in-context prompting (indicating human preference) become positive in fine-tuning (indicating AI preference). The writer type \times reader type interactions in the in-context prompting models reveal that lay readers are substantially more favorable to AI-generated text than experts, a difference that disappears after fine-tuning.

S6.2: Primary Hypothesis Tests

Table 9 presents the preregistered contrasts testing H1 and H2, corresponding to panels A-B in Figure 2. The odds ratios reveal a striking reversal: expert readers' 6-8 fold preference for human writing in in-context prompting conditions (OR = 0.16 for style, 0.13 for quality) transforms into an 8-fold preference for AI in fine-tuning (OR = 8.16 for style). Lay readers show less dramatic but directionally consistent shifts.

S6.3: Model-Predicted Probabilities

Table 10 reports the predicted probabilities displayed in Figure 2 panels C-D. These probabilities, derived from the inverse-logit transformation of linear predictors, show the convergence of expert and lay preferences after fine-tuning. While experts strongly prefer humans over all baseline models in in-context prompting settings (probabilities 0.17-0.43), fine-tuning elevates AI preference to approximately 0.74 for both reader groups.

Table 8: GLM coefficients with CR2 robust standard errors for all primary models. Negative coefficients indicate preference for the reference category (Human), while positive coefficients indicate preference for AI. The sign reversal between in-context prompting and fine-tuned models represents the core finding.

Term	Est. (log-odds)	SE	<i>z</i>	<i>p</i>
<i>Style — In-context Prompting (n = 2,400)</i>				
(Intercept)	0.901	0.148	6.098	1.07e-09
writer_typeGPT4o_baseline	-1.655	0.402	-4.120	3.78e-05
writer_typeClaude_baseline	-1.390	0.276	-5.039	4.67e-07
writer_typeGemini_baseline	-2.510	0.426	-5.894	3.77e-09
reader_typeLay	-0.826	0.169	-4.877	1.08e-06
writer_typeGPT4o_baseline:reader_typeLay	1.468	0.446	3.292	9.95e-04
writer_typeClaude_baseline:reader_typeLay	1.428	0.336	4.249	2.15e-05
writer_typeGemini_baseline:reader_typeLay	2.211	0.475	4.650	3.32e-06
<i>Style — Fine-tuned (n = 1,440)</i>				
(Intercept)	-1.050	0.141	-7.435	1.04e-13
writer_typeGPT4o_fine-tuned	2.100	0.282	7.435	1.04e-13
reader_typeLay	-0.008	0.186	-0.042	0.967
writer_typeGPT4o_fine-tuned:reader_typeLay	0.015	0.372	0.042	0.967
<i>Quality — In-context Prompting (n = 2,400)</i>				
(Intercept)	0.978	0.174	5.634	1.76e-08
writer_typeGPT4o_baseline	-2.406	0.473	-5.083	3.71e-07
writer_typeClaude_baseline	-1.246	0.368	-3.390	6.997e-04
writer_typeGemini_baseline	-2.406	0.421	-5.711	1.13e-08
reader_typeLay	-1.197	0.196	-6.124	9.12e-10
writer_typeGPT4o_baseline:reader_typeLay	3.065	0.518	5.919	3.23e-09
writer_typeClaude_baseline:reader_typeLay	1.723	0.420	4.105	4.04e-05
writer_typeGemini_baseline:reader_typeLay	2.594	0.473	5.480	4.25e-08
<i>Quality — Fine-tuned (n = 1,440)</i>				
(Intercept)	-0.314	0.122	-2.571	0.0102
writer_typeGPT4o_fine-tuned	0.627	0.244	2.571	0.0102
reader_typeLay	-0.129	0.157	-0.821	0.412
writer_typeGPT4o_fine-tuned:reader_typeLay	0.258	0.314	0.821	0.412

Table 9: Primary contrasts testing H1 (baseline LLMs vs Human) and H2 (fine-tuned vs Human). Holm correction applied within each outcome and hypothesis across reader types. The OR column shows odds ratios with values < 1 favoring humans and > 1 favoring AI.

Outcome	Hyp.	Reader	Est.	SE	<i>p</i>	<i>p</i> _{Holm}	OR	95% CI
style	H1	Expert	-1.852	0.308	3.16e-09	1.26e-08	0.157	[0.085, 0.290]
style	H1	Lay	-0.150	0.166	0.365	0.365	0.861	[0.621, 1.193]
style	H2	Expert	2.101	0.276	2.19e-14	2.19e-14	8.163	[4.693, 14.198]
style	H2	Lay	2.116	0.241	1.76e-18	1.76e-18	8.290	[5.155, 13.333]
quality	H1	Expert	-2.021	0.405	1.16e-07	2.33e-07	0.133	[0.063, 0.280]
quality	H1	Lay	0.441	0.180	0.014	0.019	1.554	[1.092, 2.212]
quality	H2	Expert	0.626	0.248	0.012	0.016	1.873	[1.161, 3.021]
quality	H2	Lay	0.886	0.197	7.08e-06	1.42e-05	2.424	[1.644, 3.574]

Table 10: Predicted probabilities of selecting AI excerpts by model and reader type, corresponding to Figure 2C-D. Values above 0.5 indicate AI preference. Note the convergence of expert and lay preferences in fine-tuned models.

Outcome	Setting	Model	Reader	\hat{p}	95% CI
style	In_Context	GPT-4o (In-Context)	Expert	0.320	[0.214, 0.449]
style	In_Context	GPT-4o (In-Context)	Lay	0.472	[0.407, 0.538]
style	In_Context	Claude (In-Context)	Expert	0.380	[0.303, 0.464]
style	In_Context	Claude (In-Context)	Lay	0.528	[0.464, 0.591]
style	In_Context	Gemini (In-Context)	Expert	0.167	[0.097, 0.271]
style	In_Context	Gemini (In-Context)	Lay	0.444	[0.374, 0.516]
style	Fine_tuned	GPT-4o (Fine-tuned)	Expert	0.741	[0.684, 0.790]
style	Fine_tuned	GPT-4o (Fine-tuned)	Lay	0.742	[0.694, 0.785]
quality	In_Context	GPT-4o (In-Context)	Expert	0.193	[0.112, 0.313]
quality	In_Context	GPT-4o (In-Context)	Lay	0.608	[0.541, 0.671]
quality	In_Context	Claude (In-Context)	Expert	0.433	[0.329, 0.544]
quality	In_Context	Claude (In-Context)	Lay	0.564	[0.499, 0.627]
quality	In_Context	Gemini (In-Context)	Expert	0.193	[0.124, 0.289]
quality	In_Context	Gemini (In-Context)	Lay	0.492	[0.422, 0.563]
quality	Fine_tuned	GPT-4o (Fine-tuned)	Expert	0.578	[0.519, 0.635]
quality	Fine_tuned	GPT-4o (Fine-tuned)	Lay	0.609	[0.562, 0.654]

S6.4: Interaction Diagnostics

Tables 11 and 12 examine the writer type \times reader type interactions in detail. The significant interactions in in-context prompting models (all $p < 0.001$) quantify lay readers' greater tolerance for AI-generated text. The absence of significant interactions in fine-tuned models ($p > 0.4$) indicates that fine-tuning produces text that both expert and lay readers find equally compelling.

Table 11: Individual interaction terms showing differential preferences between expert and lay readers. Large positive coefficients in in-context prompting models indicate lay readers are more favorable to AI than experts.

Outcome	Setting	Term	Est.	SE	z	p
style	In_Context	GPT4o_baseline:reader_typeLay	1.468	0.446	3.292	9.95e-04
style	In_Context	Claude_baseline:reader_typeLay	1.428	0.336	4.249	2.15e-05
style	In_Context	Gemini_baseline:reader_typeLay	2.211	0.475	4.650	3.32e-06
style	Fine_tuned	GPT4o_fine-tuned:reader_typeLay	0.015	0.372	0.042	0.967
quality	In_Context	GPT4o_baseline:reader_typeLay	3.065	0.518	5.919	3.23e-09
quality	In_Context	Claude_baseline:reader_typeLay	1.723	0.420	4.105	4.04e-05
quality	In_Context	Gemini_baseline:reader_typeLay	2.594	0.473	5.480	4.25e-08
quality	Fine_tuned	GPT4o_fine-tuned:reader_typeLay	0.258	0.314	0.821	0.412

The results in this section provide compelling statistical evidence for the paper's central finding: fine-tuning on author-specific corpora fundamentally transforms AI-generated text from clearly inferior (as judged by experts) to preferred over human writing. The convergence of expert and lay preferences after fine-tuning suggests that the improvements are not merely superficial but represent genuine advances in literary quality and stylistic fidelity.

S7: Author-Level Heterogeneity (Figure 3)

This section examines variation in AI preference rates across the 30 fine-tuned authors. Despite substantial differences in training corpus sizes (0.89M to 10.9M tokens) and fine-tuning costs (\$22 to \$273), we find no detectable relationship between data quantity and model performance within this token range, suggesting that even authors with limited published works can be effectively emulated.

Table 12: Joint Wald tests for writer type \times reader type interactions. Significant interactions in in-context prompting models become non-significant after fine-tuning, indicating convergence of expert and lay preferences.

Outcome	Setting	Term	df	χ^2	p
style	In_Context	writer_type:reader_type	3	24.923	1.60e-05
style	Fine_tuned	writer_type:reader_type	1	0.0017	0.967
quality	In_Context	writer_type:reader_type	3	37.584	3.46e-08
quality	Fine_tuned	writer_type:reader_type	1	0.674	0.412

S7.1: Per-Author Success Rates

Table 13 presents author-level AI preference rates using Jeffreys prior estimates to stabilize small-sample authors. The results reveal striking heterogeneity. For **style**, AI preference rates range from 18% (Tony Tulathimutte) to 98% (Roxane Gay), with 27 of 30 authors showing AI superiority (rate > 0.5). For **quality**, the range spans 30% (Ian McEwan) to 86% (Cheryl Strayed), with 23 of 30 authors favoring AI.

Notably, some authors show divergent patterns across outcomes. Tony Tulathimutte represents an extreme case: readers found his *style* uniquely difficult to emulate (18% AI preference) yet rated AI *quality* as acceptable (54%). This suggests certain idiosyncratic voices resist algorithmic mimicry even when technical writing competence is achieved.

Table 13: Per-author AI preference rates (Jeffreys-smoothed and raw) ranked by performance. Values above 0.5 indicate AI preference. Note the wide variation across authors and the divergence between **style** and quality rankings for some authors.

Outcome	Author	Rank	AI Win Rate (Jeffreys)	Human Win Rate (Jeffreys)	AI Win Rate (Raw)
quality	Cheryl Strayed	1.0	0.86	0.14	0.875
quality	Marilynne Robinson	2.0	0.78	0.22	0.792
quality	Colson Whitehead	3.0	0.74	0.26	0.750
quality	Han Kang	4.0	0.74	0.26	0.750
quality	Haruki Murakami	5.0	0.74	0.26	0.750
quality	Junot Diaz	6.0	0.74	0.26	0.750
quality	Rachel Cusk	7.0	0.74	0.26	0.750
quality	Salman Rushdie	8.0	0.74	0.26	0.750
quality	Sigrid Nunez	9.0	0.74	0.26	0.750
quality	Orhan Pamuk	10.0	0.70	0.30	0.708
quality	Lydia Davis	11.0	0.66	0.34	0.667
quality	Percival Everett	12.0	0.66	0.34	0.667
quality	Jonathan Franzen	13.0	0.62	0.38	0.625
quality	Louise Erdrich	14.0	0.62	0.38	0.625
quality	Annie Proulx	15.0	0.58	0.42	0.583
quality	George Saunders	16.0	0.58	0.42	0.583
quality	Zadie Smith	17.0	0.58	0.42	0.583
quality	Chimamanda Ngozi Adichie	18.0	0.54	0.46	0.542
quality	Jhumpa Lahiri	19.0	0.54	0.46	0.542
quality	Margaret Atwood	20.0	0.54	0.46	0.542
quality	Roxane Gay	21.0	0.54	0.46	0.542
quality	Sally Rooney	22.0	0.54	0.46	0.542
quality	Tony Tulathimutte	23.0	0.54	0.46	0.542
quality	Min Jin Lee	24.0	0.46	0.54	0.458
quality	Annie Ernaux	25.0	0.42	0.58	0.417
quality	Ben Lerner	26.0	0.42	0.58	0.417
quality	Kazuo Ishiguro	27.0	0.42	0.58	0.417
quality	Otessa Moshfegh	28.0	0.38	0.62	0.375
quality	Yoko Ogawa	29.0	0.34	0.66	0.333
quality	Ian McEwan	30.0	0.30	0.70	0.292
style	Roxane Gay	1.0	0.98	0.02	1.000
style	Chimamanda Ngozi Adichie	2.0	0.94	0.06	0.958
style	Han Kang	3.0	0.94	0.06	0.958
style	Margaret Atwood	4.0	0.94	0.06	0.958
style	Ben Lerner	5.0	0.90	0.10	0.917
style	Junot Diaz	6.0	0.90	0.10	0.917
style	Marilynne Robinson	7.0	0.90	0.10	0.917
style	Kazuo Ishiguro	8.0	0.86	0.14	0.875
style	Lydia Davis	9.0	0.86	0.14	0.875
style	Orhan Pamuk	10.0	0.82	0.18	0.833
style	Cheryl Strayed	11.0	0.78	0.22	0.792
style	Min Jin Lee	12.0	0.78	0.22	0.792
style	Sigrid Nunez	13.0	0.78	0.22	0.792

Continued on next page

Outcome	Author	Rank	AI Win Rate (Jeffreys)	Human Win Rate (Jeffreys)	AI Win Rate (Raw)
style	Colson Whitehead	14.0	0.74	0.26	0.750
style	Haruki Murakami	15.0	0.74	0.26	0.750
style	Ian McEwan	16.0	0.74	0.26	0.750
style	Rachel Cusk	17.0	0.74	0.26	0.750
style	Sally Rooney	18.0	0.74	0.26	0.750
style	Percival Everett	19.0	0.70	0.30	0.708
style	Jhumpa Lahiri	20.0	0.66	0.34	0.667
style	Otessa Moshfegh	21.0	0.66	0.34	0.667
style	Zadie Smith	22.0	0.66	0.34	0.667
style	Annie Proulx	23.0	0.62	0.38	0.625
style	George Saunders	24.0	0.62	0.38	0.625
style	Jonathan Franzen	25.0	0.62	0.38	0.625
style	Salman Rushdie	26.0	0.62	0.38	0.625
style	Yoko Ogawa	27.0	0.62	0.38	0.625
style	Louise Erdrich	28.0	0.50	0.50	0.500
style	Annie Ernaux	29.0	0.42	0.58	0.417
style	Tony Tulathimutte	30.0	0.18	0.82	0.167

Note. Jeffreys prior adds 0.5 to successes and 0.5 to failures; differences from raw are small but avoid 0/1 edge cases.

S7.2: Relationship with Corpus Size

Table 14 presents OLS regression results examining the relationship between fine-tuning corpus size (ranging from 0.89M tokens for Tulathimutte to 10.9M for Pamuk) and AI preference rates. The near-zero slopes with confidence intervals spanning zero and minimal R^2 values indicate no detectable relationship within this token range between training data quantity and model performance.

Table 14: OLS regression of author-level AI win rate on fine-tuning corpus size (millions of tokens). Slopes are near zero with CIs spanning zero and very low R^2 , indicating no detectable relationship between corpus size and AI preference rate.

Outcome	Slope	Slope CI (Low)	Slope CI (High)	Intercept	Intercept CI (Low)	Intercept CI (High)	R^2	N Authors
quality	0.0070	-0.0174	0.0314	0.570	0.458	0.681	0.0122	30
style	0.0032	-0.0262	0.0326	0.729	0.595	0.863	0.0017	30

Non-parametric check. Spearman $\rho = 0.138$ ($p = 0.467$) for quality and $\rho = -0.140$ ($p = 0.462$) for style.

This finding has economic implications: authors with limited published works (costing \$22–\$50 to fine-tune) can be emulated as effectively as prolific authors (costing \$200+). The absence of any detectable corpus-size effect ($R^2 \approx 0$) suggests that capturing an author’s voice depends more on stylistic consistency than sheer data volume.

S8: Detection Mechanisms and Economic Analysis (Figure 4)

This section examines the mechanisms linking AI detectability to human preferences (H3) and quantifies the economic implications of fine-tuning. We show that fine-tuning reverses and substantially attenuates the negative association between AI detection and preference observed in in-context prompting, while achieving performance that readers prefer at $\approx 99.7\%$ lower cost than human writers.

S8.1: AI Detection and Preference Relationship

Table 15 presents the full model examining how AI detectability (Pangram score) influences preferences. The key finding is the significant interaction between *pangram_score* and *setting* (quality: $\hat{\beta} \approx 2.90$, $p < 0.001$), indicating that fine-tuning reverses and substantially attenuates the negative association between detectability and preference observed in in-context prompting.

Table 15: Pangram GLM coefficients with CR2 robust standard errors (clustered by reader). Negative coefficients for *pangram_score* in in-context prompting indicate detection reduces preference; positive interactions with *setting* show that fine-tuning reverses and attenuates this penalty.

Outcome	N	N Readers	Term	Estimate	SE	z	p	OR	OR (Low)	OR (High)
quality	3840	159	(Intercept)	0.990	0.169	5.87	0.000	2.69	1.93	3.74
quality	3840	159	<i>pangram_score</i>	-2.010	0.333	-6.03	0.000	0.13	0.07	0.26
quality	3840	159	<i>setting</i> (Fine-tuned)	-1.020	0.175	-5.86	0.000	0.36	0.26	0.51
quality	3840	159	<i>reader_type</i> (Lay)	-1.220	0.190	-6.42	0.000	0.30	0.20	0.43
quality	3840	159	<i>pangram_score</i> × <i>setting</i>	2.900	0.876	3.32	0.001	18.20	3.28	102.00
quality	3840	159	<i>pangram_score</i> × <i>reader_type</i>	2.480	0.377	6.56	0.000	11.90	5.67	24.90
quality	3840	159	<i>setting</i> × <i>reader_type</i>	1.250	0.198	6.33	0.000	3.50	2.37	5.15
quality	3840	159	<i>pangram_score</i> × <i>setting</i> × <i>reader_type</i>	-3.290	1.030	-3.21	0.001	0.04	0.00	0.28
style	3840	159	(Intercept)	0.909	0.146	6.23	0.000	2.48	1.86	3.31
style	3840	159	<i>pangram_score</i>	-1.850	0.292	-6.32	0.000	0.16	0.09	0.28
style	3840	159	<i>setting</i> (Fine-tuned)	-0.938	0.158	-5.95	0.000	0.39	0.29	0.53
style	3840	159	<i>reader_type</i> (Lay)	-0.832	0.167	-4.97	0.000	0.44	0.31	0.60
style	3840	159	<i>pangram_score</i> × <i>setting</i>	2.560	0.810	3.15	0.002	12.90	2.63	63.00
style	3840	159	<i>pangram_score</i> × <i>reader_type</i>	1.690	0.336	5.03	0.000	5.41	2.80	10.50
style	3840	159	<i>setting</i> × <i>reader_type</i>	0.840	0.178	4.73	0.000	2.32	1.64	3.28
style	3840	159	<i>pangram_score</i> × <i>setting</i> × <i>reader_type</i>	-1.900	0.922	-2.06	0.039	0.15	0.02	0.91

Notes. CR2 robust standard errors clustered by reader. Readers: 28 Expert and 131 Lay (total $N_{\text{readers}} = 159$).

S8.2: Predicted Probabilities Across Detection Levels

Tables 16 and 17 show model-predicted probabilities at different AI detection levels, corresponding to Figure 2F. For expert readers evaluating *quality*, high detection scores (0.9) reduce AI preference to ~31% in in-context prompting but have a point estimate of ~68% after fine-tuning (95% CI 0.37–0.89), demonstrating that fine-tuning makes outputs robust to detection-based skepticism.

Table 16: Predicted AI preference probabilities for *style* at different Pangram detection scores. Fine-tuning *reverses and attenuates* the relationship with detection score (with wider uncertainty at high detection).

Bin	Setting	Reader Type	Prob	Prob (Low)	Prob (High)
0.1	In-context Prompting	Expert	0.674	0.621	0.722
0.5	In-context Prompting	Expert	0.497	0.494	0.500
0.9	In-context Prompting	Expert	0.320	0.273	0.372
0.1	Fine-tuned	Expert	0.511	0.491	0.530
0.5	Fine-tuned	Expert	0.581	0.435	0.714
0.9	Fine-tuned	Expert	0.648	0.380	0.847
0.1	In-context Prompting	Lay	0.515	0.483	0.547
0.5	In-context Prompting	Lay	0.500	0.499	0.500
0.9	In-context Prompting	Lay	0.484	0.451	0.517
0.1	Fine-tuned	Lay	0.507	0.495	0.520
0.5	Fine-tuned	Lay	0.557	0.461	0.649
0.9	Fine-tuned	Lay	0.606	0.428	0.759

Table 17: Predicted AI preference probabilities for *quality*. Expert readers show strong detection sensitivity in in-context prompting (~0.69 → ~0.31) that *reverses* after fine-tuning (point estimate at 0.9 ~0.68; 95% CI 0.37–0.89).

Bin	Setting	Reader Type	Prob	Prob (Low)	Prob (High)
0.1	In-context Prompting	Expert	0.688	0.628	0.742
0.5	In-context Prompting	Expert	0.496	0.493	0.499
0.9	In-context Prompting	Expert	0.306	0.254	0.363
0.1	Fine-tuned	Expert	0.514	0.490	0.537
0.5	Fine-tuned	Expert	0.602	0.431	0.751
0.9	Fine-tuned	Expert	0.683	0.373	0.887
0.1	In-context Prompting	Lay	0.454	0.421	0.488
0.5	In-context Prompting	Lay	0.501	0.500	0.501
0.9	In-context Prompting	Lay	0.547	0.512	0.582
0.1	Fine-tuned	Lay	0.501	0.487	0.515
0.5	Fine-tuned	Lay	0.509	0.403	0.614
0.9	Fine-tuned	Lay	0.517	0.324	0.705

S8.3: Stylometric Correlates and Mediation

Table 18 shows that cliché density has the strongest correlation with AI detection (Pearson $r = 0.60$), while readability ease is negatively correlated ($r = -0.23$), suggesting AI text is marked by formulaic phrases but simpler syntax. In in-context prompting conditions, approximately 16% of the detection effect on preference is mediated through cliché density; after fine-tuning, this mediation drops to a statistically insignificant 1.3%, indicating that fine-tuning substantially reduces rather than merely masking these stylistic signatures.

Table 18: Correlations between Pangram AI detection scores and stylometric features.

Metric	Pearson r	Spearman ρ	N
readability_ease	-0.232	-0.312	330
cliche_density	0.600	0.602	330
total_adjective_count	0.104	0.129	330
num_cliches	0.614	0.633	330

S8.4: Economic Analysis

Table 19 documents the cost structure for AI-based text generation. With fine-tuning costs ranging from \$22.25 to \$272.50 (median \$77.88) plus \$3 for inference to generate 100,000 words, the total AI cost represents approximately 0.3% of the \$25,000 a professional writer would charge. This ~99.7% cost reduction, combined with reader preference for fine-tuned outputs, quantifies the potential economic disruption to creative writing markets.

Table 19: Cost comparison for generating 100,000 words. Fine-tuning plus inference costs (\$25–\$276) represent less than 1% of professional writer compensation (\$25,000).

N Authors	Fine-tune Min	Fine-tune Med	Fine-tune Max	Total Min	Total Med	Total Max
30	\$22.25	\$77.88	\$272.50	\$25.25	\$80.88	\$275.50

Note. Inference cost of \$3 per 100k words follows `fig4c_cost_summary.csv`; per-token assumptions are documented in the repository.

S9. Deviations from Preregistration

This section documents deviations from the preregistered analysis plan (<https://osf.io/zt4ad>). All other aspects were implemented exactly as specified.

- **Model framework.** The preregistration specified mixed-effects logistic regression with random intercepts for readers and prompts. We implemented logistic GLMs with CR2 cluster-robust standard errors (clustered by reader) due to convergence issues and singularity warnings in the mixed-effects models. Point estimates from both approaches were similar when convergence was achieved, and inference on the preregistered contrasts (H1, H2) remained unchanged.
- **Sample sizes (exceeded targets).** We successfully recruited more participants than planned:
 - Expert readers: 28 (planned: 25)
 - Lay readers: 131 (planned: 120)
 - Fine-tuned authors: 30 (planned minimum: 10)

All analyses used the full realized sample. No stopping rules were applied or violated; the additional data strengthens the reliability of our findings.

- **Additional analyses (not preregistered).** We added:
 - Writer type \times Reader type interaction tests (Type-III Wald with robust covariance) for transparency
 - Predicted probability visualizations (Figures 2C-D) to complement the odds ratio panels

These additions provide fuller context but do not alter the conclusions drawn from the preregistered H1 and H2 contrasts.

S10. Code and Data Availability

All code and data necessary to reproduce Figures 2-4 and the associated statistical analyses will be made publicly available on paper acceptance. Please contact authors if you need it:

- **Preregistration:** <https://osf.io/zt4ad> (Version 2, July 2025)
- **Key Analysis Scripts:** The repository contains R scripts for data processing, model fitting, and figure generation. Core analyses can be reproduced by running the numbered scripts in sequence.
- **Environment:** Analyses were conducted in R 4.3.1 with key packages including `clubSandwich` (CR2 robust SEs) and `emmeans` (contrasts). Full package versions are documented in the repository.
- **Data:** Anonymized trial-level data in long format, with documentation of all variable definitions and transformations applied.